

# Camlp5 - Reference Manual

version 4.06

Daniel de Rauglaudre  
July 27, 2007

Copyright © 2007 Institut National de Recherche en Informatique et Automatique  
This document was generated by a shell script from the html documentation pages.



# Contents

<b>1</b>	<b>Overview</b>	<b>5</b>
1.1	Shell usage . . . . .	5
1.2	Parsing and Printing kits . . . . .	6
1.3	Extending syntax . . . . .	6
1.4	Pretty printing . . . . .	6
1.5	Note: the revised syntax . . . . .	7
<b>2</b>	<b>Parsing and Printing tools</b>	<b>9</b>
2.1	Stream parsers . . . . .	9
2.2	Extensible grammars . . . . .	10
2.3	Pretty module . . . . .	10
<b>3</b>	<b>Syntax extensions</b>	<b>13</b>
3.1	Entries . . . . .	13
3.2	Syntax tree quotations . . . . .	14
3.3	An example . . . . .	14
3.3.1	The code . . . . .	15
3.3.2	Compilation . . . . .	16
3.3.3	Testing . . . . .	16
<b>4</b>	<b>The revised syntax</b>	<b>19</b>
4.1	Modules, Structure and Signature items . . . . .	19
4.2	Expressions and Patterns . . . . .	20
4.2.1	Imperative constructions . . . . .	20
4.2.2	Tuples and Lists . . . . .	20
4.2.3	Records . . . . .	21
4.2.4	Irrefutable patterns . . . . .	21
4.2.5	Constructions with matching . . . . .	21
4.2.6	Mutables and Assignment . . . . .	22
4.2.7	Miscellaneous . . . . .	22
4.3	Types and Constructors . . . . .	23
4.4	Streams and Parsers . . . . .	24
4.5	Classes and Objects . . . . .	25
4.6	Labels and Variants . . . . .	26
<b>5</b>	<b>Quotations</b>	<b>27</b>
5.1	Introduction . . . . .	27
5.2	Quotation expander . . . . .	27
5.3	Defining a quotation . . . . .	28
5.3.1	By syntax tree . . . . .	28

5.3.2	By string . . . . .	28
5.3.3	Default quotation . . . . .	28
5.4	Antiquotations . . . . .	29
5.4.1	Example without antiquotation node . . . . .	29
5.4.2	Example with antiquotation node . . . . .	30
5.4.3	In conclusion . . . . .	31
5.5	A full example: lambda terms . . . . .	31
5.5.1	Lexer . . . . .	32
5.5.2	Parser . . . . .	33
5.5.3	Compilation and test . . . . .	34
<b>6</b>	<b>Stream parsers</b>	<b>37</b>
6.1	Introduction . . . . .	37
6.2	Syntax . . . . .	37
6.3	Streams . . . . .	38
6.3.1	Stream.from . . . . .	38
6.3.2	Stream.of_list . . . . .	38
6.3.3	Stream.of_string . . . . .	38
6.3.4	Stream.of_channel . . . . .	38
6.4	Semantics of parsers . . . . .	39
6.4.1	Parser . . . . .	39
6.4.2	Match with parser . . . . .	40
6.4.3	Error messages . . . . .	40
6.4.4	Stream pattern component . . . . .	41
6.4.5	Let statement . . . . .	42
6.4.6	Lookahead . . . . .	42
6.4.7	No error optimization . . . . .	43
6.4.8	Position . . . . .	44
6.4.9	Semantic action . . . . .	44
6.5	Remarks . . . . .	44
6.5.1	Simplicity vs Associativity . . . . .	44
6.5.2	Lexing vs Parsing . . . . .	45
6.5.3	Lexer syntax vs Parser syntax . . . . .	46
6.5.4	Purely functional parsers . . . . .	46
<b>7</b>	<b>Functional parsers</b>	<b>47</b>
7.1	Syntax . . . . .	47
7.2	Streams . . . . .	48
7.2.1	Fstream.from . . . . .	48
7.2.2	Fstream.of_list . . . . .	48
7.2.3	Fstream.of_string . . . . .	48
7.2.4	Fstream.of_channel . . . . .	48
7.3	Semantics of parsers . . . . .	48
7.3.1	Fparser . . . . .	48
7.3.2	Error position . . . . .	48
<b>8</b>	<b>Stream lexers</b>	<b>51</b>
8.1	Introduction . . . . .	51
8.2	Syntax . . . . .	51
8.3	Semantics . . . . .	52
8.3.1	Symbols . . . . .	53
8.3.2	Specific expressions . . . . .	54

8.3.3	Lookahead . . . . .	54
8.3.4	Semantic actions of rules . . . . .	55
8.3.5	A complete example . . . . .	55
8.3.6	Compiling . . . . .	55
8.3.7	How to display the generated code . . . . .	56
<b>9</b>	<b>Extensible grammars</b>	<b>57</b>
9.1	Getting started . . . . .	57
9.2	Syntax of the EXTEND statement . . . . .	58
9.3	Semantics of the EXTEND statement . . . . .	59
9.3.1	GLOBAL indicator . . . . .	59
9.3.2	Entries list . . . . .	60
9.3.3	Rules insertion . . . . .	62
9.3.4	Semantic action . . . . .	62
9.4	The DELETE_RULE statement . . . . .	63
9.5	Extensions FOLD0 and FOLD1 . . . . .	63
9.6	Extensions SLIST0, SLIST1 and SOPT . . . . .	63
9.7	Grammar machinery . . . . .	63
9.7.1	Start and Continue . . . . .	63
9.7.2	Associativity . . . . .	64
9.7.3	Errors and recovery . . . . .	64
9.7.4	Tokens starting rules . . . . .	65
9.7.5	Kind of grammar . . . . .	65
9.8	The Grammar module . . . . .	65
9.8.1	Main types and values . . . . .	65
9.8.2	Printing grammar entries . . . . .	66
9.8.3	Clearing grammars and entries . . . . .	66
9.8.4	Functorial interface . . . . .	67
9.8.5	Miscellaneous . . . . .	68
9.9	Interface with the lexer . . . . .	68
9.9.1	Token patterns . . . . .	69
9.9.2	The glexer record . . . . .	69
9.9.3	Minimalist version . . . . .	71
9.10	Functorial interface . . . . .	71
9.10.1	The glexer type . . . . .	72
9.10.2	The functor parameter . . . . .	72
9.10.3	The resulting grammar module . . . . .	72
9.10.4	GEXTEND and GDELETE_RULE . . . . .	73
<b>10</b>	<b>Pretty module</b>	<b>75</b>
10.1	Module description . . . . .	75
10.1.1	horiz_vertic . . . . .	75
10.1.2	sprintf . . . . .	75
10.1.3	line_length . . . . .	76
10.2	Example . . . . .	76
10.3	Programming with Pretty . . . . .	76
10.3.1	Hints . . . . .	76
10.3.2	How to cancel a horizontal print . . . . .	77
10.4	Remarks . . . . .	77
10.4.1	Kernel . . . . .	77
10.4.2	Strings vs Channels . . . . .	77
10.4.3	Strings or other types . . . . .	78

10.4.4	Why raising exceptions ? . . . . .	78
<b>11</b>	<b>Printing programs</b>	<b>79</b>
11.1	Introduction . . . . .	79
11.2	Principles . . . . .	79
11.2.1	Using module Pretty . . . . .	79
11.2.2	Using module Extfun and its syntax . . . . .	81
11.2.3	Dangling else, bar, semicolon . . . . .	82
11.2.4	By level . . . . .	82
<b>12</b>	<b>Scheme and Lisp syntaxes</b>	<b>85</b>
12.1	Common . . . . .	85
12.2	Scheme syntax . . . . .	86
12.3	Lisp syntax . . . . .	86
<b>13</b>	<b>Syntax tree</b>	<b>87</b>
13.1	Introduction . . . . .	87
13.2	Location . . . . .	88
13.2.1	In expressions . . . . .	88
13.2.2	In patterns . . . . .	88
13.3	Antiquotations . . . . .	89
13.4	Nodes and Quotations . . . . .	89
13.4.1	expr . . . . .	90
13.4.2	patt . . . . .	91
13.4.3	ctyp . . . . .	92
13.4.4	modules... . . . .	92
13.4.5	classes... . . . .	94
13.4.6	other . . . . .	95
<b>14</b>	<b>Locations</b>	<b>97</b>
14.1	Definitions . . . . .	97
14.2	Building locations . . . . .	97
14.3	Raising with a location . . . . .	98
14.4	Other functions . . . . .	98
<b>15</b>	<b>Pragma directive</b>	<b>101</b>
<b>16</b>	<b>Extensible functions</b>	<b>103</b>
16.1	Syntax . . . . .	103
16.2	Semantics . . . . .	103

# Chapter 1

## Overview

Camlp5 is a preprocessor and pretty-printer for ocaml programs.

As a preprocessor, it allows to:

- add syntax extensions,
- redefine the whole syntax of the language.

As a pretty printer, it allows to:

- display programs in an elegant way,
- convert from a syntax to another,
- check the results of syntax extensions.

Camlp5 also provides some parsing and pretty printing tools.

It works as a shell command and can also be used in the ocaml toplevel.

### 1.1 Shell usage

The main shell commands are:

- `camlp5o` : to treat files written in normal ocaml syntax,
- `camlp5r` : to treat files written in an original syntax named the *revised syntax*.

These commands can be given as parameters of the option `-pp` of the ocaml compiler. Examples:

```
ocamlc -pp camlp5o foo.ml
ocamlc -pp camlp5r bar.ml
```

This way, the parsing is done by camlp5. In case of syntax errors, the parsing fails with an error message and the compilation is aborted. Otherwise, the ocaml compiler continues with the syntax tree provided by camlp5.

In the toplevel, it is possible to preprocess the input phrases by loading one of the files `"camlp5o.cma"` or `"camlp5r.cma"`. The common usage is:

```
ocaml -I +camlp5 camlp5o.cma
ocaml -I +camlp5 camlp5r.cma
```

## 1.2 Parsing and Printing kits

Parsing and printing extensions are ocaml object files, i.e. files with the extension ".cmo" or ".cma". They are the result of the compilation of ocaml source files containing what is necessary to do the parsing or printing extensions. These object files are named *parsing* or *printing kits*.

These files cannot be linked to produce executables because they generally call functions and use variables defined only in camlp5 core, typically belonging to the module "Pcaml". The kits are destined to be loaded by the camlp5 commands, either by their command arguments or by the source files using them through the directive "#load".

It is therefore important to compile the *kits* with the option "-c" of the ocaml compiler (i.e. just compilation, do not produce executable) and with the option "-I +camlp5" to inform the compiler to find module interfaces in installed camlp5 library.

In the ocaml toplevel, it is possible to use a kit by simply loading it with the directive "#load".

## 1.3 Extending syntax

A syntax extension is a camlp5 parsing kit. There are two ways to use a syntax extension:

- Either by giving this object file as parameter to the camlp5 command. For example:

```
ocamlc -pp "camlp5o ./myext.cmo" foo.ml
```

- Or by adding the directive "#load" in the source file:

```
#load ". /myext.cmo";;
```

and compile it simply like this:

```
ocamlc -pp camlp5o foo.ml
```

Several syntax extensions can be used for a single file. The way to create one's own syntax extensions is explained further.

## 1.4 Pretty printing

Like for syntax extensions, the pretty printing are defined or extended through camlp5 printing kits. The main pretty printing kits provided by camlp5 are:

- `pr_o.cmo`: to pretty print in normal syntax,
- `pr_r.cmo`: to pretty print in revised syntax.

Examples: if we have a file, `foo.ml`, written in normal syntax and another one, `bar.ml`, written in revised syntax, here are the commands to pretty print them in their own syntax:

```
camlp5o pr_o.cmo foo.ml  
camlp5r pr_r.cmo bar.ml
```

And how to convert them into the other syntax:

```
camlp5o pr_r.cmo foo.ml  
camlp5r pr_o.cmo foo.ml
```

The way to create one's own pretty printing extensions is explained further.

## 1.5 Note: the revised syntax

The *revised syntax* is a specific syntax whose aim is to resolve some syntax problems and inconsistencies of the normal ocaml syntax. A chapter will explain the differences between the normal and the revised syntax.

All examples of this documentation are written in revised syntax. Even if you don't know it, it is not difficult to understand. The same examples can be written in normal syntax. In case of problems, refer to the chapter describing it.



# Chapter 2

## Parsing and Printing tools

Camlp5 provides two original parsing tools:

- stream parsers
- extensible grammars

The first parsing tool, the stream parsers, is the elementary system. It is pure syntactic sugar, i.e. the code is directly converted into basic ocaml statements: functions, pattern matchings, try. A stream parser is a function. But the system does not take care of associativity, nor parsing level, and left recursion result on infinite loops, just like functions whose first action would be a call to itself.

The second parsing tool, the extensible grammars, are more sophisticated. A grammar written with them is more readable, and look like grammars written with tools like "yacc". They take care of associativity, left recursion, and level of parsing. They are dynamically extensible, what allows the syntax extensions what camlp5 provides for ocaml syntax.

In both cases, the input data are streams.

Camlp5 also provides a pretty printing tool, a module allowing to control the lines length.

The next sections give an overview of the two parsing and the pretty tools.

### 2.1 Stream parsers

The stream parsers are a system of elementary recursive descendant parsing. Streams are actually lazy lists. At each step, the head of the list is compared against a *stream pattern*. There are two kinds of streams parsers:

- The imperative streams parsers, where the elements are removed from the stream as long as they are parsed. Parsers return either:
  - A value, in case of success,
  - The exception "**Stream.Failure**" when the parser does not apply and no elements have been removed from the stream, indicating that, possibly, other parsers may apply,
  - The exception "**Stream.Error**" when the parser does not apply, but one or several elements have been removed from the stream, indicating that nothing can be done to make up the error.

- The purely functional stream parsers where the elements are not removed from the stream during the parsing. These parsers return a value of type "option", i.e either:
  - "Some" a value and the remaining stream, in case of success,
  - "None", in case of failure.

The differences are about:

- Syntax errors: in the imperative version, the location of the error is clear, it is at the current position of the stream, and the system allows to provide a specific error message (typically, that some "element" was "expected"). On the other hand, in the functional version, the position is not clear since it returns with nothing more than the initial stream. The only solution to know where the error happened is to analyze that stream to see how many elements have be unfrozen, and no clear error message is available, just "syntax error".
- Power: in the imperative version, when a rule raises the exception "**Stream.Error**", the parsing cannot continue. In the functional version, the parsing can continue by analyzing the next rule with the initial unaffected stream: this is a *limited backtrack*.
- Neatness: functional streams are neater, just like functional programming is neater than imperative programming.

In the imperative version, there exists also lexers, a shorter syntax when the stream elements are of the specific type 'char'.

## 2.2 Extensible grammars

Extensible grammars manipulate *grammar entries*. Grammar entries are abstract values internally containing mutable stream parsers. When a grammar entry is created, its internal parser is empty, i.e. it raises "**Parse.Failure**" if used. A specific syntactic construction, with the keyword "**EXTEND**" allow to extend grammar entries with new grammar rules.

In opposition to stream parsers, grammar entries take care of associativity, left factorization, and levels. Moreover, the syntax for grammars allows to define optional calls, lists and lists with separators. However, they are not functions and cannot take parameters.

Since the internal system is stream parsers, extensible grammars use recursive descendant parsing.

The parser of the ocaml language in `camlp5` is written with extensible grammars.

## 2.3 Pretty module

The "**Pretty**" module is an original tool allowing to control the displaying of lines. The user has to specify two functions where:

- the data is printed in one only line
- the data is printed in several lines

The system first tries the first function. At any time, it the line overflows, i.e. if its size is greater than some "line length" specified in the module interface, or if it contains newlines, the function is aborted and control is given to the second function.

This is a basic, but powerful, system. It supposes that the programmer takes care of the current indentation, and the beginning and the end of its lines.

The module will be extended in the future to hide the management of indentations and line continuations, and by the supply of functions combining the two cases above, in which the programmer can specify the possible places where newlines can be inserted.



# Chapter 3

## Syntax extensions

Camlp5 allows to extend the syntax of the ocaml language, and even change the whole syntax.

It uses for that one of its parsing tools: the extensible grammars.

To understand the whole syntax in the examples given in this chapter, it is a good thing to know this parsing tool, but we shall try to give some minimal explanations to allow the reader to follow them.

A syntax extension is an ocaml object file (ending with ".cmo" or ".cma") which is loaded inside camlp5. The source of this file uses calls to the specific statement `EXTEND` applied to entries defined in the camlp5 module `"Pcaml"`.

### 3.1 Entries

The grammar of ocaml contains several entries, corresponding to the major notions of the language, which are modifiable this way, and even erasable. They are defined in this module `"Pcaml"`.

Most important entries:

- `expr`: the expressions.
- `patt`: the patterns.
- `ctyp`: the types.
- `str_item`: the structure items, i.e. the items between "struct" and "end", and the toplevel phrases in a ".ml" file.
- `sig_item`: the signature items, i.e. the items between "sig" and "end", and the toplevel phrases in a ".mli" file.
- `module_expr`: the module expressions.
- `module_type`: the module types.

Entries of object programming:

- `class_expr`: the class expressions.
- `class_type`: the class types.

- `class_str_item`: the objects items.
- `class_sig_item`: the objects types items.

Main entries of files and interactive toplevel parsing:

- `implem`: the phrases that can be found in a ".ml" file.
- `interf`: the phrases that can be found in a ".mli" file.
- `top_phrase`: the phrases of the interactive toplevel.
- `use_file`: the phrases that can be found in a file included by the directive "use".

Extra useful entries also accessible:

- `let_binding`: the bindings "expression = pattern" found in the "let" statement.
- `type_declaration`: the bindings "name = type" found in the "type" statement.

## 3.2 Syntax tree quotations

A grammar rule is a list of rule symbols followed by the semantic action, i.e. the result of the rule. This result is a syntax tree, whose type is the type of the extended entry. The description of the types of the syntax tree are in the `camlp5` module "MLast".

However, there is a simpler way to make values of these syntax tree types: the system quotations. With this system, it is possible to represent syntax tree in concrete syntax, between specific parentheses.

For example, the syntax node of the "if" statement, supposing you know the meaning of the constructors of the "MLast" module would be:

```
MLast.ExIf loc e1 e2 e3
```

where `loc` is the source location, and `e1`, `e2`, `e3` are the expressions constituting the if statement. With quotations, we could write it like this:

```
<:expr< if $e1$ then $e2$ else $e3$ >>
```

With quotations, it is possible to build pieces of program as complex as desired. See the chapter about syntax trees.

## 3.3 An example

A classical extension is the programming of the "repeat" statement. The "repeat" statement is like a "while" except that the loop is executed at least one time and that the test is at the end of the loop and is inverted. The equivalent of:

```
repeat x; y; z until c
```

is:

```
do {
  x; y; z;
  while not c do { x; y; z }
}
```

or, with a loop:

```
loop () where rec loop () = do {  
  x; y; z;  
  if c then () else loop ()  
}
```

### 3.3.1 The code

This syntax extension could be written like this (see the detail of syntax in the chapter about extensible grammars and the syntax tree quotations in the chapter about them):

```
#load "pa_extend.cmo";  
#load "q_MLast.cmo";  
open Pcaml;  
EXTEND  
  expr:  
    [ [ "repeat"; el = LIST1 expr SEP ";"; "until"; c = expr ->  
      let el = el @ [<:expr< while not $c$ do { $list:el$ } >>] in  
      <:expr< do { $list:el$ } >> ] ]  
  ;  
END;
```

Alternatively, with the loop version:

```
#load "pa_extend.cmo";  
#load "q_MLast.cmo";  
open Pcaml;  
EXTEND  
  expr:  
    [ [ "repeat"; el = LIST1 expr SEP ";"; "until"; c = expr ->  
      let el = el @ [<:expr< if $c$ then () else loop () >>] in  
      <:expr< loop () where rec loop () = do { $list:el$ } >> ] ]  
  ;  
END;
```

The first "#load" in the code means that a syntax extension has been used in the file, namely the "EXTEND" statement. The second "#load" means that abstract tree quotations has been used, namely the "<:expr< ... >>".

The quotation in the second version is especially interesting. If we wanted to write it with the abstract syntax tree, we would have to write it like this:

```
MLast.ExLet loc True  
  [(MLast.PaLid loc "loop",  
    MLast.ExFun loc [(MLast.PaUId loc "()", None, MLast.ExSeq loc el)])]  
  (MLast.ExApp loc (MLast.ExLid loc "loop") (MLast.ExUId loc "()));
```

which is less easy for the programmer and less readable.

### 3.3.2 Compilation

If the file "foo.ml" contains one of these versions, it is possible to compile it like this:

```
ocamlc -pp camlp5r -I +camlp5 -c foo.ml
```

Notice that the ocamlc option "-c" means that we are interested only in generating the object file "foo.cmo", not achieving the compilation which, anyway, would not link because of usage of modules specific to camlp5.

### 3.3.3 Testing

In the ocaml toplevel

```
ocaml -I +camlp5 camlp5r.cma
      Objective Caml version ...

      Camlp5 Parsing version ...

# #load "foo.cmo";
# value x = ref 42;
value x : ref int = {val=42}
# repeat
  print_int x.val; print_endline ""; x.val := x.val + 3
until x.val > 70;
42
45
48
51
54
57
60
63
66
69
- : unit = ()
```

In a file

The code, above, used in the toplevel, can be written in a file, say "bar.ml":

```
#load "./foo.cmo";
value x = ref 42;
repeat
  print_int x.val;
  print_endline "";
  x.val := x.val + 3
until x.val > 70;
```

with a subtle difference: the loaded file must be "./foo.cmo" and not just "foo.cmo" because camlp5 does not have, by default, the current directory in its path.

The file can be compiled like this:

```
ocamlc -pp camlp5r bar.ml
```

or in native code:

```
ocamlopt -pp camlp5r bar.ml
```

... work in progress ...



# Chapter 4

## The revised syntax

The revised syntax is an alternative syntax of ocaml. It is close to the normal syntax. We present here only the differences between the two syntaxes.

Notice that there is a simple way to know how the normal syntax is written in revised syntax: write the code in a file "foo.ml" in normal syntax and type, in a shell:

```
camlp5o pr_r.cmo pr_rp.cmo foo.ml
```

And, conversely, how a file "bar.ml" written in revised syntax is displayed in normal syntax:

```
camlp5r pr_o.cmo pr_op.cmo bar.ml
```

Even simpler, without creating a file:

```
camlp5o pr_r.cmo pr_op.cmo -impl -  
... type in normal syntax ...  
... type control-C ...  
camlp5r pr_o.cmo pr_rp.cmo -impl -  
... type in revised syntax ...  
... type control-C ...
```

### 4.1 Modules, Structure and Signature items

- Structure and signature items always end with a single semicolon which is required.
- In structures, the declaration of a value is introduced by the keyword "value", instead of "let":

OCaml	Revised
<code>let x = 42;;</code> <code>let x = 42 in x + 7;;</code>	<code>value x = 42;</code> <code>let x = 42 in x + 7;</code>

- In signatures, the declaration of a value is also introduced by the keyword "value", instead of "val":

OCaml	Revised
<code>val x : int;;</code>	<code>value x : int;</code>

- In signatures, abstract module types are represented by a quote and an (any) identifier:

OCaml	Revised
<code>module type MT;;</code>	<code>module type MT = 'a;</code>

- Functor application uses curryfication. Parentheses are not required for the parameters:

OCaml	Revised
<code>type t = Set.Make(M).t;; module M = Mod.Make (M1) (M2);;</code>	<code>type t = (Set.Make M).t; module M = Mod.Make M1 M2;</code>

- It is possible to group several declarations together either in an interface or in an implementation by enclosing them between "declare" and "end" (this is useful when using syntax extensions to generate several declarations from one). Example in an interface:

```
declare
  type foo = [ Foo of int | Bar ];
  value f : foo -> int;
end;
```

## 4.2 Expressions and Patterns

### 4.2.1 Imperative constructions

- The sequence is introduced by the keyword "do" followed by "{" and terminated by "}"; it is possible to put a semicolon after the last expression:

OCaml	Revised
<code>e1; e2; e3; e4</code>	<code>do { e1; e2; e3; e4 }</code>

- The "do" after the "while" loop and the "for" loop are followed by a "{" and the loop end with "}"; it is possible to put a semicolon after the last expression:

OCaml	Revised
<code>while e1 do   e1; e2; e3 done  for i = e1 to e2 do   e1; e2; e3 done</code>	<code>while e1 do {   e1; e2; e3 }  for i = e1 to e2 do {   e1; e2; e3 }</code>

### 4.2.2 Tuples and Lists

- Parentheses are required in tuples:

OCaml	Revised
<code>1, "hello", World</code>	<code>(1, "hello", World)</code>

- The lists are always enclosed with brackets. A list is a left bracket, followed by a list of elements separated with semicolons, optionally followed by colon-colon and an element, and ended by a right bracket. Warning: the colon-colon is not an infix but is just part of the syntactic construction.

OCaml	Revised
<code>x :: y [x; y; z] x :: y :: z :: t</code>	<code>[x :: y] [x; y; z] [x; y; z :: t]</code>

### 4.2.3 Records

- In record update, parentheses are required around the initial expression:

OCaml	Revised
<code>{e with field = a}</code>	<code>{(e) with field = a}</code>

- It is authorized to use function binding syntax in record field definitions:

OCaml	Revised
<code>{field = fun a -&gt; e}</code>	<code>{field a = e}</code>

### 4.2.4 Irrefutable patterns

An *irrefutable pattern* is a pattern where it is syntactically visible that it never fails. They are used in some syntactic constructions. It is either:

- A variable,
- The wildcard `"_"`,
- The constructor `"()"`,
- A tuple with irrefutable patterns,
- A record with irrefutable patterns,
- A type constraint with an irrefutable pattern.

Notice that this definition is only syntactic: a constructor belonging to a type having only one constructor is not considered as an irrefutable pattern (except `"()"`).

### 4.2.5 Constructions with matching

- The keyword `"function"` no longer exists. Only `"fun"` is used.
- The pattern matchings, in constructions with `"fun"`, `"match"` and `"try"` are closed with brackets: an open bracket `"["` before the first case, and a close bracket `"]"` after the last case:

OCaml	Revised
<code>match e with   p1 -&gt; e1   p2 -&gt; e2</code>	<code>match e with [ p1 -&gt; e1   p2 -&gt; e2 ]</code>

But if there is only one case and if the pattern is irrefutable, the brackets are not required. These examples work identically in OCaml and in revised syntax:

OCaml	Revised
<code>fun x -&gt; x fun {foo=(y, _)} -&gt; y</code>	<code>fun x -&gt; x fun {foo=(y, _)} -&gt; y</code>

- It is possible to write the empty function, raising the exception `"Match_failure"` whichever parameter is applied, the empty `"match"`, raising `"Match_failure"` after having evaluated its expression, and the empty `"try"`, equivalent to its expression without `try`:

```

fun []
match e with []
try e with []

```

- The patterns after "let" and "value" must be irrefutable. The following OCaml expression:

```
let f (x::y) = ...
```

must be written:

```
let f = fun [ [x::y] -> ...
```

- It is possible to use a construction "where", equivalent to "let", but usable only when where is only one binding. The expression:

```
e1 where p = e
```

is equivalent to:

```
let p = e in e1
```

## 4.2.6 Mutables and Assignment

- The statement "<-" is written ":=":

OCaml	Revised
<code>x.f &lt;- y</code>	<code>x.f := y</code>

- The "ref" type is declared as a record type with one field named "val", instead of "contents". The operator "!" does not exist any more, and references are assigned like the other mutables:

OCaml	Revised
<code>x := !x + y</code>	<code>x.val := x.val + y</code>

## 4.2.7 Miscellaneous

- The "else" is required in the "if" statement:

OCaml	Revised
<code>if a then b</code>	<code>if a then b else ()</code>

- The boolean operations "or" and "and" can only be written with "||" and "&&":

OCaml	Revised
<code>a or b &amp; c</code>	<code>a    b &amp;&amp; c</code>
<code>a    b &amp;&amp; c</code>	<code>a    b &amp;&amp; c</code>

- No more "begin end" construction. One must use parentheses.
- The operators as values are written with an backslash:

OCaml	Revised
<code>(+)</code>	<code>\+</code>
<code>(mod)</code>	<code>\mod</code>

- Nested "as" patterns require parenthesis:

OCaml	Revised
<code>function Some a as b, c -&gt;</code> ...	<code>fun [ ((Some a as b), c) -&gt;</code> ...

But they are not required before the right arrow:

OCaml	Revised
<code>function Some a as b -&gt;</code> ...	<code>fun [ Some a as b -&gt;</code> ...

- The operators with special characters are not automatically infix. To define infixes, use the syntax extensions.

## 4.3 Types and Constructors

- The type constructors are before their type parameters, which are curryfied:

OCaml	Revised
<code>int list</code> <code>('a, bool) Hashtbl.t</code> <code>type 'a foo = 'a list list</code>	<code>list int</code> <code>Hashtbl.t 'a bool</code> <code>type foo 'a = list (list 'a)</code>

- The abstract types are represented by a unbound type variable:

OCaml	Revised
<code>type 'a foo;;</code> <code>type bar;;</code>	<code>type foo 'a = 'b;</code> <code>type bar = 'a;</code>

- Parentheses are required in tuples of types:

OCaml	Revised
<code>int * bool</code>	<code>(int * bool)</code>

- In declaration of a concrete type, brackets must enclose the constructor declarations:

OCaml	Revised
<code>type t = A of i   B;;</code>	<code>type t = [ A of i   B ];</code>

- It is possible to make the empty type, without constructor:

```
type foo = [];
```

- There is a syntax difference between data constructors with several parameters and data constructors with one parameter of type tuple:

The declaration of a data constructor with several parameters is done by separating the types with "and". In expressions and patterns, this constructor parameters must be curryfied:

OCaml	Revised
<code>type t = C of t1 * t2;;</code> <code>C (x, y);;</code>	<code>type t = [ C of t1 and t2 ];</code> <code>C x y;</code>

The declaration of a data constructor with one parameter of type tuple is done by using a tuple type. In expressions and patterns, the parameter has not to be curried, since it is alone. In that case the syntax of constructor parameters is the same between the two syntaxes:

OCaml	Revised
<code>type t = D of (t1 * t2);; D (x, y);;</code>	<code>type t = [ D of (t1 * t2) ]; D (x, y);</code>

- The bool constructors start with an uppercase letter. The identifiers "true" and "false" are not keywords:

OCaml	Revised
<code>true &amp;&amp; false</code>	<code>True &amp;&amp; False</code>

- In record types, the keyword "mutable" must appear after the colon:

OCaml	Revised
<code>type t = {mutable x : t1};;</code>	<code>type t = {x : mutable t1};</code>

- Manifest types are with "==":

OCaml	Revised
<code>type 'a t = 'a option = None   Some of 'a</code>	<code>type t 'a = option 'a == [ None   Some of 'a ]</code>

- Polymorph types start with "!":

OCaml	Revised
<code>type t = { f : 'a . 'a list }</code>	<code>type t = { f : ! 'a . list 'a }</code>

## 4.4 Streams and Parsers

- The streams and the stream patterns are bracketed with "[:" and ":]" instead of "<" and ">".
- The stream component "terminal" is written with a backquote instead of a quote:

OCaml	Revised
<code>[&lt; '1; '2; s; '3 &gt;]</code>	<code>[[: '1; '2; s; '3 :]</code>

- The cases of parsers are bracketed with "[" and "]", like for "fun", "match" and "try". If there is one case, the brackets are not required:

OCaml	Revised
<code>parser [&lt; 'Foo &gt;] -&gt; e   [&lt; p = f &gt;] -&gt; f;; parser [&lt; 'x &gt;] -&gt; x;;</code>	<code>parser [ [: 'Foo :] -&gt; e   [: p = f :] -&gt; f ]; parser [: 'x :] -&gt; x;</code>

- It is possible to write the empty parser raising the exception "Stream.Failure" whichever parameter is applied, and the empty stream matching always raising "Stream.Failure":

```

parser []
match e with parser []

```

- In normal syntax, the error indicator starts with a double question mark, in revised syntax with a simple question mark:

OCaml	Revised
<pre> parser   [&lt; '1; '2 ?? "error" &gt;] -&gt;   ... </pre>	<pre> parser   [: '1; '2 ? "error" :] -&gt;   ... </pre>

- In normal syntax, the component optimization starts with "?!", in revised syntax with "!":

OCaml	Revised
<pre> parser   [&lt; '1; '2 ?! &gt;] -&gt;   ... </pre>	<pre> parser   [: '1; '2 ! :] -&gt;   ... </pre>

## 4.5 Classes and Objects

- Object items end with a single semicolon which is required.
- Class type parameters follow the class identifier:

OCaml	Revised
<pre> class ['a, 'b] point = ... class c = [int] color;; </pre>	<pre> class point ['a, 'b] = ... class c = color [int]; </pre>

- In the type of class with parameters, the type of the parameters are between brackets. Example in signature:

OCaml	Revised
<pre> class c : int -&gt; point;; </pre>	<pre> class c : [int] -&gt; point; </pre>

- The keywords "virtual" and "private" must be in this order:

OCaml	Revised
<pre> method virtual private m :   ... method private virtual m :   ... </pre>	<pre> method virtual private m :   ... method virtual private m :   ... </pre>

- Object variables are introduced with "value" instead of "val":

OCaml	Revised
<pre> object val x = 3 end </pre>	<pre> object value x = 3; end </pre>

- Type constraints in objects are introduced with "type" instead of "constraint":

OCaml	Revised
<pre> object constraint 'a = int end </pre>	<pre> object type 'a = int; end </pre>

## 4.6 Labels and Variants

- Labels in types must start with " ":

OCaml	Revised
<code>val x : num:int -&gt; bool;;</code>	<code>value x : num:int -&gt; bool;</code>

- Types whose number of variants are fixed start with "[ =":

OCaml	Revised
<code>type t = ['On   'Off];;</code>	<code>type t = [ = 'On   'Off];</code>

- The "[]" and the "<" in variant types must not be stucked:

OCaml	Revised
<code>type t = [&lt; 'Foo   'Bar ];;</code>	<code>type t = [ &lt; 'Foo   'Bar ];</code>

# Chapter 5

## Quotations

The quotations are a syntax extension in camlp5 allowing to build expressions and patterns in any syntax independant from the one of ocaml. Quotations are *expanded*, i.e. transformed, at parse time to produce normal syntax trees, like the rest of the program. Quotations *expanders* are normal ocaml functions writable by any programmer.

The aim of quotations is to use concrete syntax for manipulating abstract values, what make programs easier to write, read, modify, and understand. Their drawback is that they are isolated from the rest of the program, in opposition to syntax extensions, which are included in the language.

### 5.1 Introduction

A quotation is syntactically enclosed by specific quotes formed by less (<) and greater (>) signs. Namely:

- starting with either "<<" or "<:ident<" where "ident" is the quotation name,
- ending with ">>"

Examples:

```
<< \x.x x >>  
<:foo< hello, world >>  
<:bar< @#$%;* >>
```

The text between these particular parentheses can be any text. It may contain enclosing quotations and the characters "<", ">" and "\" can be escaped by "\". Notice that possible double-quote, parentheses, ocaml comments do not have necessary to balance inside them.

As far as the lexer is concerned, a quotation is just a kind of string.

### 5.2 Quotation expander

The quotations are treated at parse time. Each quotation name is associated with a *quotation expander*, a function transforming the content of the quotation into a syntax tree. There are actually two expanding functions, depending on the fact that the quotation is in the context of an expression or if it is in the context of a pattern.

If a quotation has no associated quotation expander, a parsing error is displayed and the compilation fails.

The quotation expander, or, rather, expanders, are functions taking the quotation string as parameter and returning a syntax tree. There is no constraint about which parsing technology is used. It can be stream parsers, extensible grammars, string scanning, ocamllex and yacc, any.

To build syntax trees, camlp5 provides a way to easily build them: see the chapter about them: it is possible to build abstract syntax through concrete syntax using, precisely... quotations.

## 5.3 Defining a quotation

### 5.3.1 By syntax tree

To define a quotation, it is necessary to program the quotation expanders and to, finally, end the source code with a call to:

```
Quotation.add name (Quotation.ExAst (f_expr, f_patt));
```

where "name" is the name of the quotation, and "f\_expr" and "f\_patt" the respective quotations expanders for expressions and patterns.

Then, after compilation of the source file (without linking, i.e. using option "-c" of the ocaml compiler), an object file is created (ending with ".cmo"), which can be used as syntax extension *kit* of camlp5.

### 5.3.2 By string

There is, actually, another way to program the expander: an alone function which returns, not a syntax tree, but just a string which is parsed, afterwards, by the system. This function takes a boolean as first parameter telling whether the quotation is in position of expression (True) or in position of a pattern (False).

If using that way, the source file must end with:

```
Quotation.add name (Quotation.ExStr f);
```

where "f" is that quotation expander. The advantage of this second approach is that it is simple to understand and use. The drawback is that there is no way to specify different source locations for different parts of the quotation (what may be important in semantic error messages).

### 5.3.3 Default quotation

It is possible to use some quotation without its name. Use for that the variable "Quotation.default\_quotation". For example, ending a file by:

```
Quotation.add "foo" (Quotation.ExAst (f_expr, f_patt));
Quotation.default.val := "foo";
```

allows to use the quotation "foo" without its name, i.e.:

```
<< ... >>
```

instead of:

```
<:foo< ... >>
```

If several files set the variable "Quotation.default", the default quotation applies to the last loaded one.

## 5.4 Antiquotations

A quotation obeys its own rules of lexing and parsing. Its result is a syntax tree, of type `Pcaml.expr` if the quotation is in the context of an expression, or `Pcaml.patt` if the quotation is in the context of a pattern.

But it can be interesting to insert portions of expressions or patterns of the enclosing context in its own quotations. For that, the syntax of the quotation must define a syntax for *antiquotations areas*. It can be, for example:

- A character introducing a variable: in this case the antiquotation can just be a variable.
- A couple of characters enclosing the antiquotations. For example, in the predefined syntax tree quotations, the antiquotations are enclosed with dollar (`"$"`) signs.

In quotations, the locations in the resulting syntax tree are all set to the location of the quotation itself (if this resulting tree contains locations, they are overwritten with this location). Consequently, if there are semantic (typing) errors, the ocaml compiler will underline the entire quotation.

But in antiquotations, since they can be inserted in the resulting syntax tree, it is interesting to keep their initial quotations. For that, the nodes:

```
<:expr< $anti:...$ >>  
<:patt< $anti:...$ >>
```

equivalent to:

```
MLast.ExAnt loc ...  
MLast.PaAnt loc ...
```

are provided (see syntax tree quotations).

Let us take an example, without this node, and with this specific node.

Let us consider an elementary quotation system whose contents is just an antiquotation. This is just a school example, since the quotations brackets are not necessary, in this case. But we are going to see how the source code is underlined in errors messages.

### 5.4.1 Example without antiquotation node

The code for this quotation is (file `"qa.ml"`):

```
#load "q_MLast.cmo";  
let expr s = Grammar.Entry.parse Pcaml.expr (Stream.of_string s) in  
Quotation.add "a" (Quotation.ExAst (expr, fun []));
```

The quotation expander `"expr"` just takes the string parameter (the contents of the quotation), and returns the result of the expression parser of the ocaml language.

Compilation:

```
ocamlc -pp camlp5r -I +camlp5 -c qa.ml
```

Let us test it in the toplevel with a voluntary typing error:

```
$ ocaml -I +camlp5 camlp5r.cma
Objective Caml version ...

Camlp5 Parsing version ...

# #load "qa.cmo";
# let x = "abc" and y = 25 in <:a< x ^ y >>;
Characters 28-41:
  let x = "abc" and y = 25 in <:a< x ^ y >>;
                                ~~~~~
This expression has type int but is here used with type string
```

We observe that the full quotation is underlined, although it concerns only the variable "y".

## 5.4.2 Example with antiquotation node

Let us consider this second version (file "qb.ml"):

```
#load "q_MLast.cmo";
let expr s =
  let ast = Grammar.Entry.parse Pcaml.expr (Stream.of_string s) in
  let loc = Stdpp.make_lined_loc 1 0 (0, String.length s) in
  <:expr< $anti:ast$ >>
in
Quotation.add "b" (Quotation.ExAst (expr, fun []));
```

Like above, the quotation expander "expr" takes the string parameter (the contents of the quotation) and applies the expression parser of the ocaml language. But its result, instead of being returned as it is, is enclosed with the antiquotation node. (The location built is the location of the whole string.)

Compilation:

```
ocamlc -pp camlp5r -I +camlp5 -c qb.ml
```

Now the same test gives:

```
$ ocaml -I +camlp5 camlp5r.cma
Objective Caml version ...

Camlp5 Parsing version ...

# #load "qb.cmo";
# let x = "abc" and y = 25 in <:b< x ^ y >>;
Characters 37-38:
  let x = "abc" and y = 25 in <:b< x ^ y >>;
                                ~
This expression has type int but is here used with type string
```

Notice that, now, only the variable "y" is underlined.

### 5.4.3 In conclusion

In the resulting tree of the quotation expander:

- only portions of this tree, which are sons of the `expr` and `patt` antiquotation nodes, have the right location they have in the quotation (provided the quotation expander gives it the right location of the antiquotation in the quotation),
- the other nodes inherit, as location, the location of the full quotation.

## 5.5 A full example: lambda terms

This example allows to represent lambda terms by a concrete syntax and to be able to combine them using antiquotations.

A lambda term is defined like this:

```
type term =  
  [ Lam of string and term  
  | App of term and term  
  | Var of string ]  
;
```

Examples:

```
value fst = Lam "x" (Lam "y" (Var "x"));  
value snd = Lam "x" (Lam "y" (Var "y"));  
value delta = Lam "x" (App (Var "x") (Var "x"));  
value omega = App delta delta;  
value comb_s =  
  Lam "x"  
    (Lam "y"  
      (Lam "z"  
        (App (App (Var "x") (Var "y")) (App (Var "x") (Var "z"))))));
```

Since combinations of lambda term may be complicated, The idea is to represent them by quotations in concrete syntax. We want to be able to write the examples above like this:

```
value fst = << \x.\y.x >>;  
value snd = << \x.\y.y >>;  
value delta = << \x.x x >>  
value omega = << ^delta ^delta >>;  
value comb_s = << \x.\y.\z.(x y)(x z) >>;
```

which is a classic representation of lambda terms.

Notice, in the definition of "`omega`", the usage of the caret ("`^`") sign to specify antiquotations. Notice also the simplicity of the representation of the expression defining "`comb_s`".

Here is the code of the quotation expander, `term.ml`. The expander uses the extensible grammars. It has its own lexer (using the stream lexers) because the lexer of ocaml programs ("`Plexer.gmake ()`"), cannot recognize the backslashes alone.

### 5.5.1 Lexer

```
(* lexer *)

#load "pa_lex.cmo";

value rev_implode l =
  let s = String.create (List.length l) in
  loop (String.length s - 1) l where rec loop i =
    fun
      [ [c :: l] -> do { String.unsafe_set s i c; loop (i - 1) l }
      | [] -> s ]
  ;

module B =
  struct
    value empty = [];
    value add x l = [x :: l];
    value get = rev_implode;
  end
;

value rec ident =
  lexer
  [ "a..zA..Z0..9-_'\128..\255" ident! | ]
;

value empty _ = parser [: _ = Stream.empty :] -> [];

value rec next_tok =
  lexer
  [ "\\\" -> ("BSLASH", "")
  | "^" -> ("CARET", "")
  | "a..z" ident! -> ("IDENT", $buf)
  | "(" -> ("(", "(")
  | ")" -> ("", ")")
  | "." -> ("", ".")
  | empty -> ("EOS", "")
  | -> raise (Stream.Error "lexing error: bad character") ]
;

value rec skip_spaces = lexer [ " \n\r"/ skip_spaces! | ];

value record_loc loct i (bp, ep) = do {
  if i >= Array.length loct.val then do {
    let newt =
      Array.init (2 * Array.length loct.val + 1)
      (fun i ->
        if i < Array.length loct.val then loct.val.(i)
        else Stdpp.dummy_loc)
    in
    loct.val := newt;
  }
}
```

```

    }
    else ();
    loct.val.(i) := Stdpp.make_loc (bp, ep)
};

value lex_func cs =
  let loct = ref [| |] in
  let ts =
    Stream.from
      (fun i -> do {
        ignore (skip_spaces $empty cs : list char);
        let bp = Stream.count cs in
        let r = next_tok $empty cs in
        let ep = Stream.count cs in
        record_loc loct i (bp, ep);
        Some r
      })
  in
  (ts, fun i -> loct.val.(i))
;

value lam_lex =
  {Token.tok_func = lex_func;
   Token.tok_using _ = (); Token.tok_removing _ = ();
   Token.tok_match = Token.default_match;
   Token.tok_text = Token.lexer_text;
   Token.tok_comm = None}
;

```

### 5.5.2 Parser

```

(* parser *)

#load "pa_extend.cmo";
#load "q_MLast.cmo";

value g = Grammar.gcreate lam_lex;
value expr_term_eos = Grammar.Entry.create g "term";
value patt_term_eos = Grammar.Entry.create g "term";

EXTEND
  GLOBAL: expr_term_eos patt_term_eos;
  expr_term_eos:
    [ [ x = expr_term; EOS -> x ] ]
  ;
  expr_term:
    [ [ BSLASH; i = IDENT; "."; t = SELF -> <:expr< Lam $str:i$ $t$ >> ]
    | [ x = SELF; y = SELF -> <:expr< App $x$ $y$ >> ]
    | [ i = IDENT -> <:expr< Var $str:i$ >>
      | CARET; r = expr_antiquot -> r
      | "("; t = SELF; ")" -> t ] ]
  ;

```

```

expr_antiquot:
  [ [ i = IDENT ->
      let r =
        let loc = Stdpp.make_loc (0, String.length i) in
        <:expr< $lid:i$ >>
      in
        <:expr< $anti:r$ >> ] ]
;
patt_term_eos:
  [ [ x = patt_term; EOS -> x ] ]
;
patt_term:
  [ [ BSLASH; i = IDENT; "."; t = SELF -> <:patt< Lam $str:i$ $t$ >> ]
  | [ x = SELF; y = SELF -> <:patt< App $x$ $y$ >> ]
  | [ i = IDENT -> <:patt< Var $str:i$ >>
      | CARET; r = patt_antiquot -> r
      | "("; t = SELF; ")" -> t ] ]
;
patt_antiquot:
  [ [ i = IDENT ->
      let r =
        let loc = Stdpp.make_loc (0, String.length i) in
        <:patt< $lid:i$ >>
      in
        <:patt< $anti:r$ >> ] ]
;
END;

value expand_expr s = Grammar.Entry.parse expr_term_eos (Stream.of_string s);
value expand_patt s = Grammar.Entry.parse patt_term_eos (Stream.of_string s);

Quotation.add "term" (Quotation.ExAst (expand_expr, expand_patt));
Quotation.default.val := "term";

```

### 5.5.3 Compilation and test

Compilation:

```
ocamlc -pp camlp5r -I +camlp5 -c term.ml
```

Example, in the toplevel, including a semantic error, correctly underlined, thanks to the antiquotation nodes:

```

$ ocaml -I +camlp5 camlp5r.cma
Objective Caml version ...

Camlp5 Parsing version ...

# #load "term.cmo";
# type term =
  [ Lam of string and term
  | App of term and term
  | Var of string ]

```

```

;
type term = [ Lam of string and term | App of term and term | Var of string ]
# value comb_s = << \x.\y.\z.(x y)(x z) >>;
value comb_s : term =
  Lam "x"
    (Lam "y"
      (Lam "z" (App (App (Var "x") (Var "y")) (App (Var "x") (Var "z"))))))
# value omega = << ^delta ^delta >>;
Characters 18-23:
  value omega = << ^delta ^delta >>;
                ^^^^^
Unbound value delta
# value delta = << \x.x x >>;
value delta : term = Lam "x" (App (Var "x") (Var "x"))
# value omega = << ^delta ^delta >>;
value omega : term =
  App (Lam "x" (App (Var "x") (Var "x"))) (Lam "x" (App (Var "x") (Var "x")))

```



# Chapter 6

## Stream parsers

We describe here the syntax and the semantics of the parsers of streams of Camlp5. Streams are kinds of lazy lists. The parsers of these streams use recursive descent method, which is the most natural one in functional languages. In particular, parsers are normal functions.

Notice that the parsers have existed in OCaml since many years (the beginning of the 90ies), but some new features have been added in 2007 (lookahead, "no error" optimization, and let..in statement) in Camlp5 distribution. This page describes them also.

### 6.1 Introduction

Parsers apply to values of type "Stream.t" defined in the module "Stream" of the standard library of OCaml. Like the type "list", the type "Stream.t" has a type parameter, indicating the type of its elements. They differ from the lists that they are lazy (the elements are evaluated as long as the parser need them for its actions), and imperative (parsers deletes their first elements when they take their parsing decisions): notice that pure functional parsers exist in Camlp5, where the corresponding streams are lazy and functional, the analyzed elements remaining in the initial stream and the semantic action returning the resulting stream together with the normal result, which allow natural limited backtrack but have the drawback that it is not easy to find the position of parsing errors when they happen.

Parsers of lazy/imperative streams, which are described here, use a method named "recursive descent": they look at the first element, they decide what to do in function of its value, and continue the parsing with the remaining elements. Parsers can call other parsers, and can be recursive, like normal functions.

Actually, parsers are just pure syntactic sugar. When writing a parser in the syntax of the parser, Camlp5 transforms them into normal call to functions, use of patterns matchings and try..with statements. The pretty printer of Camlp5, by default, display this expanded result, without syntax of parsers. A pretty printing kit, when added, can rebuild the parsers in their initial syntax and display it.

### 6.2 Syntax

The syntax of the parsers, when loading "pa\_rp.cmo" (or already included in the command "camlp5r"), is the following:

```
expression ::= parser
            | match-with-parser
parser ::= "parser" pos-opt "[" parser-cases "]"
```

```

        | "parser" pos-opt parser-case
match-with-parser ::= "match" expression with "parser" pos-opt "[" parser-cases "]"
        | "match" expression with "parser" pos-opt parser-case
parser-cases ::= parser-cases parser-case
        | <nothing>
parser-case ::= "[:" stream-pattern ":" pos-opt "->" expression
stream-pattern ::= stream-patt-comp
        | stream-patt-comp ";" stream-patt-cont
        | "let" LIDENT "=" expression "in" stream-pattern
        | <nothing>
stream-patt-cont ::= stream-patt-comp-err
        | stream-patt-comp-err ";" stream-patt-cont
        | "let" LIDENT "=" expression "in" stream-patt-cont
stream-patt-comp-err ::= stream-patt-comp
        | stream-patt-comp "?" expression
        | stream-patt-comp "!"
stream-patt-comp ::= "\"" pattern
        | "\"" pattern "when" expression
        | "?=" lookaheads
        | pattern "=" expression
        | pattern
lookaheads ::= lookaheads "|" lookahead
        | lookahead
lookahead ::= "[" patterns "]"
patterns ::= patterns pattern
        | pattern
pos-opt ::= pattern
        | <nothing>

```

## 6.3 Streams

The parsers are functions taking as parameters streams, which are values of type `"Stream.t a"` for some type `"a"`. It is possible to build streams using the functions defined in the module `"Stream"`:

### 6.3.1 Stream.from

`"Stream.from f"` returns a stream built from the function `"f"`. To create a new stream element, the function `"f"` is called with the current stream count, starting with zero. The user function `"f"` must return either `"Some <value>"` for a value or `"None"` to specify the end of the stream.

### 6.3.2 Stream.of\_list

Return a stream built from the list in the same order.

### 6.3.3 Stream.of\_string

Return a stream of the characters of the string parameter.

### 6.3.4 Stream.of\_channel

Return a stream of the characters read from the input channel parameter.

## 6.4 Semantics of parsers

### 6.4.1 Parser

A parser, defined with the syntax "parser" above, is of type "`Stream.t a -> b`" where "a" is the type of the elements of the streams and "b" the type of the result. The parser cases are tested in the order they are defined until one of them applies. The result is the semantic action of the parser case which applies. If no parser case applies, the exception "`Stream.Failure`" is raised.

When testing a parser case, if the first stream pattern component matches, all remaining stream pattern component of the stream pattern must match also. If one does not match, the parser return with the exception "`Stream.Error`" which has a parameter of type string: by default, this string is the empty string, but if the stream pattern component which does not match is followed by a question mark and an expression, this expression is evaluated and given as parameter to "`Stream.Error`".

In short, a parser can return with three ways:

- A normal result, of type "b" for a parser of type "`Stream.t a -> b`".
- Raising the exception "`Stream.Failure`".
- Raising the exception "`Stream.Error`".

Fundamentally, the exception "`Stream.Failure`" means "this parser does not apply and no element have been removed from the initial stream". This is a normal case when parsing: the parser fails, but the parsing can continue.

Conversely, the exception "`Stream.Error`" means that "this parsers encountered a syntax error". In this case, the parsing definitively fails.

Like said above, when a parser case applies partially but not completely, the exception "`Stream.Error`" is raised. This means that if two parser cases start with the same stream pattern components, only the first one have chances to work. This is a limitation of this parsing method. To turn over this limitation, one must factorize the rules to make both work. A typical example is a language having a rule for the statement "if..then..else" and a rule for the shorter "if..then". If written that way:

```
parser
[ [: 'If; e1 = expr; 'Then; e2 = expr; 'Else; e3 = expr :] -> f e1 e2 e3
| [: 'If; e1 = expr; 'Then; e2 = expr :] -> g e1 e2 ]
```

the second pattern case cannot work.

Indeed, if the input is a stream containg the syntax with "if..then..else", the first pattern case apply and it is ok. Conversely, if the input contains just "if..then", the matching fails when arriving at the "Else" stream pattern component above. Since it is not the first stream pattern component of the parser case, the parser fails with "`Stream.Error`".

Reversing the parser cases like this does not help:

```
parser
[ [: 'If; e1 = expr; 'Then; e2 = expr :] -> g e1 e2
| [: 'If; e1 = expr; 'Then; e2 = expr; 'Else; e3 = expr :] -> f e1 e2 e3 ]
```

In this case, whatever the input, containing the "if..then" construction, or the "if..then..else.." construction, the first parser case applies. If the input is "if..then..else..", only the part "if..then.." is removed from the stream and the part "else.." remains. When a parser case is applied, like for a function, the next case will never apply.

The solution, with this parsing technology, is to do a left factorization of the parsing cases, namely programming it like this:

```
parser
  [: 'If; e1 = expr; 'Then; e2 = expr;
    a =
      parser
        [ [: 'Else; e3 = expr :] -> f e1 e2 e3
          | [: :] -> g e1 e2 ] :] -> a
```

### 6.4.2 Match with parser

The syntax "match expression with parser" allows to match a stream against a parser. It is, for "parser", the equivalent of "match expression with" for "fun". The same way we could say:

```
match expression with ...
```

could be considered as an equivalent to:

```
(fun ...) expression
```

we could consider that:

```
match expression with ...
```

is an equivalent to:

```
(parser ...) expression
```

### 6.4.3 Error messages

A "Stream.Error" exception is raised when a stream pattern component does not match and that it is not the first one of the parser case. This exception has a parameter of type string, useful to specify the error message. By default, this is the empty string. To specify an error message, add a question mark and an expression after the stream pattern component. A typical error message is "that stream pattern component expected". Example with the parser of "if..then..else.." above:

```
parser
  [: 'If; e1 = expr ? "expression expected after 'if'";
    'Then ? "'then' expected";
    e2 = expr ? "expression expected after 'then'";
    a =
      parser
        [ [: 'Else; e3 = expr ? "expression expected" :] -> f e1 e2 e3
          | [: :] -> g e1 e2 ] :] -> a
```

Notice that the expression after the question mark is evaluated only in case of syntax error. Therefore, it can be a complicated call to a complicated function without slowing down the normal parsing.

#### 6.4.4 Stream pattern component

In a stream pattern (starting with "[:" and ending with ":]"), are the stream pattern components separated with the semicolon character. There are three cases of stream pattern components with some sub-cases for some of them, and an extra syntax can be used with a "let.in" construction. The three cases are:

- A direct test of one or several stream elements (called **terminal** symbol), in three ways:
  1. The character "backquote" followed by a pattern, meaning: if the stream starts with an element which is matched by this pattern, the stream pattern component matches, and the stream element is removed from the stream.
  2. The character "backquote" followed by a pattern, the keyword "when" and an expression of type "bool", meaning: if the stream starts with an element which is matched by this pattern and then if the evaluation of the expression is "True", the stream pattern component matches, and the first element of the stream is removed.
  3. The character "question mark" followed by the character "equal" (new feature 2007) and a lookahead expression (see further), meaning: if the lookahead applies, the stream pattern component matches. The lookahead may unfreeze one or several elements on the stream, but does not remove them.
- A pattern followed by the "equal" sign and an expression of type "Stream.t x -> y" for some types "x" and "y". This expression is called a **non terminal** symbol. It means: call the expression (which is a parser) with the current stream. If this sub-parser:
  1. Returns an element, the pattern is bound to this result and the next stream pattern component is tested.
  2. Raises the exception "Stream.Failure", there are two cases:
    - if the stream pattern component is the first one of the stream case, the current parser also fails with the exception "Stream.Failure".
    - if the stream pattern component is not the first one of the stream case, the current parser fails with the exception "Stream.Error".

In this second case:

- If the stream pattern component is followed by a "question mark" and an expression (which must be of type "string"), the expression is evaluated and given as parameter of the exception.
  - If the expression is followed by an "exclamation mark" (new feature 2007), the test and conversion from "Stream.Failure" to "Stream.Error" is not done, and the parser just raises "Stream.Failure" again. This is an optimization which must be assumed by the programmer, in general when he knows that the sub-parser called never raises "Stream.Failure" (for example if the called parser ends with a parser case containing an empty stream pattern). See "no error optionization" below.
  - Otherwise the exception parameter is the empty string
- A pattern, which is bound to the current stream.

Notice that patterns are bound immediately and can be used in the next stream pattern component.

### 6.4.5 Let statement

Between stream pattern components, it is possible to use the "let..in" construction (new feature 2007). This is not considered as a real stream pattern component, in the fact that it is not tested against the exception "Stream.Failure" it may raise. It can be useful for an intermediate computation. In particular, it is used internally by the lexers (this can be seen by printing a code using these lexers into parsers by Camlp5 using the pretty printing kit "pr\_r.cmo").

Example of use, when an expression have to be used several times (in the example, "d a", which is bound to the variable "c"):

```
parser
[ : a = b;
  let c = d a in
  e =
    parser
    [ [ : f = g : ] -> h c
      | [ : : ] -> c ] : ] -> e
```

### 6.4.6 Lookahead

The lookahead feature allows to look at several terminals in the stream without removing them, in order to take decisions when more than one terminal is necessary.

For example, if parsing a the normal syntax of the OCaml language, there is a problem, in recursing descendent parsing, for the cases where to treat and differentiate the following inputs:

```
(-x+1)
(-)
```

The first case is treated in a rule, telling: "a left parenthesis, followed by an expression, and a right parenthesis". The second one is "a left parenthesis, an operator, a right parenthesis". If programming it like this (left factorizing the first parenthesis):

```
parser
[ : 'Lparen;
  e =
    parser
    [ [ : e = expr; 'Rparen : ] -> e
      | [ : 'Minus; 'Rparen : ] -> minus_op ] : ] -> e
```

this does not work if the input is "(-)" because the rule "e = expr" accepts the minus sign as expression starting, removing it from the input stream and fails as parsing error, while encountering the right parenthesis.

Conversely, writing it this way:

```
parser
[ : 'Lparen;
  e =
    parser
    [ [ : 'Minus; 'Rparen : ] -> minus_op
      | [ : e = expr; 'Rparen : ] -> e ] : ] -> e
```

does not help, because if the input is " $(-x+1)$ " the rule above starting with "`Minus`" is accepted and the exception "`Stream.Error`" is raised while encountering the variable "`x`" since a right parenthesis is expected.

In general, these cases are resolved by a left factorization of the parser cases (see the section "Semantics" above), but it is not possible in this case. The solution is to test whether the character after the minus sign is a right parenthesis:

```
parser
  [: 'Lparen;
    e =
      parser
        [ [: ?= [ _ Rparen ]; 'Minus; 'Rparen :] -> minus_op
          | [: e = expr; 'Rparen :] -> e ] :] -> e
```

It is possible to put several lists of patterns separated by a vertical bar in the lookahead construction, but with a limitation (due to the implementation): all lists of patterns must have the same number of elements.

### 6.4.7 No error optimization

The "no error optimization" is a new feature 2007. This is the fact to end a stream pattern component of kind "non-terminal" ("pattern" "equal" "expression") by the character "exclamation mark". Like said above, this inhibits the transformation of the exception "`Stream.Failure`", possibly raised by the called parser, into the exception "`Stream.Error`".

The code:

```
parser [: a = b; c = d ! :] -> e
```

is equivalent to:

```
parser [: a = b; s :] -> let c = d s in e
```

One interest of the first syntax is that it shows to readers that "`d`" is indeed a sub-parser. In the second syntax, it is called in the semantic action, which makes the parser case not so clear, as far as readability is concerned.

If the stream pattern component is the last one of the stream pattern, this allows possible tail recursion done by the OCaml compiler, in the following case:

```
parser [: a = b; c = d ! :] -> c
```

since it is equivalent (with the fact that "`c`" is at the same time the pattern of the last case and the expression of the parser case semantic action:

```
parser [: a = b; s :] -> d s
```

The call to "`d s`" can be a tail recursive call. Without the use of the "exclamation mark" in the rule, the equivalent code is:

```
parser [: a = b; s :] ->
  try d s with [ Stream.Failure -> raise (Stream.Error "")
```

which is not tail recursive (because the "try..with" construction pushes a context), preventing the compiler to optimize its code. It can have some importance when many recursive calls happen, since it can overflow the OCaml stack.

### 6.4.8 Position

The optional "pattern" before and after a stream pattern is bound to the current stream count. Indeed, streams contains internally the count of its elements. At the beginning the count is zero. When an element is removed, the count is incremented. The example:

```
parser [: a = b :] ep -> c
```

is equivalent to:

```
parser [: a = b; s :] -> let ep = Stream.count s in c
```

There is not direct syntax equivalent to the optional pattern at beginning of the stream pattern:

```
parser bp [: a = b :] -> c
```

These optional patterns allow to dispose of the stream count at the beginning and at the end of the parser case, allowing to compute locations of the rule in the source. In particular, if the stream is a stream of characters, these counts are the source location in number of characters.

### 6.4.9 Semantic action

In a parser case, after the stream pattern, there is an "arrow" and an expression, called the "semantic action". If the parser case is matched the parser returns with the expression evaluated, in a context of the possibly patterns bound in the stream pattern (position, backquote patterns, pattern equal expression, pattern).

## 6.5 Remarks

### 6.5.1 Simplicity vs Associativity

This parsing technology has the advantage to be simple to use and to understand. But it supposes to sometimes left factorize the rules. Moreover, it does not treat the associativity of operators. For example, if you try to write a parser like this (to compute arithmetic expressions):

```
value rec expr =  
  parser  
  [ [: e1 = expr; '+'; e2 = expr :] -> e1 + e2  
    | [: '('0'..'9' as c) :] -> Char.code c - Char.code '0' ]
```

this would endless loop, exactly like if you wrote a code starting like:

```
value rec expr e =  
  let e1 = expr e in  
  ...
```

A solution to that is to treat the associativity "by hand", by reading a sub-expression, then looping with a parser case parsing the operator and another sub-expression, and so on.

Another solution is to previously write parsing "combinators". Indeed, parsers being normal functions, it is possible to make a function which takes a parser as parameter and returning a parser using it. For example, left and right associativity parsing combinators:

```

value rec left_assoc op elem =
  let rec op_elem x =
    parser
    [ [: t = op; y = elem; r = op_elem (t x y) :] -> r
      | [: :] -> x ]
  in
  parser [: x = elem; r = op_elem x :] -> r
;

value rec right_assoc op elem =
  let rec op_elem x =
    parser
    [ [: t = op; y = elem; r = op_elem y :] -> t x r
      | [: :] -> x ]
  in
  parser [: x = elem; r = op_elem x :] -> r
;

```

which can be used, e.g. like this:

```

value expr =
  List.fold_right (fun op elem -> op elem)
    [left_assoc (parser [: '+' :] -> fun x y -> x +. y);
     left_assoc (parser [: '*' :] -> fun x y -> x *. y);
     right_assoc (parser [: '^' :] -> fun x y -> x ** y)]
    (parser [: '('0'..'9' as c :] -> float (Char.code c - Char.code '0'))
;

```

and tested, e.g. in the toplevel, like that:

```
expr (Stream.of_string "2^3^2");
```

The same way, it is possible to parser non-context free grammars, by programming parsers returning other parsers.

A third solution, to resolve the problem of associativity, is to use the grammars of Camlp5, which have the other advantage that they are extensible.

## 6.5.2 Lexing vs Parsing

In general, while analyzing a language, there are two levels:

- The level where the input, considered as a stream of characters, is read to make a stream of tokens (for example "words", if it is a human language, or punctuation). This level is generally called "lexing".
- The level where on read a stream of token to parse grammar rules. This level is generally called "parsing".

The "parser" construction described here can be used for both, thanks to the polymorphism of OCaml:

- The lexing level is a "parser" of streams of characters returning tokens.
- The parsing level is a "parser" of streams of tokens returning syntax trees.

By comparison, the programs "lex" and "yacc" use two different technologies. With "parser"s, it is possible to use the same one for both.

### 6.5.3 Lexer syntax vs Parser syntax

For "lexers", i.e. for the specific case of parsers when the input is a stream of characters, it is possible to use a shorter syntax. See the page on lexers. They are another syntax for this case, adapted for the specific type "`char`" and much shorter and simple to use. But they still are internally parsers of streams with the same semantics.

### 6.5.4 Purely functional parsers

This system of parsers are not purely functional in the sense that the stream structure is lazy and imperative: while parsing, the stream advances and the terminals already parsed disappear from the stream structure. This is useful because it is not necessary to return the remaining stream together with the normal result. And it is the reason why there is this "**Stream.Error**" exception: when it happens, it means that some terminals have been consumed from the stream, definitively lost, and that therefore it is no more possible to try other parser cases.

An alternative of that is using purely functional parsers using a new stream type, lazy but not destructive. Their advantage is that they use a limited backtrack: the case of "if..then..else.." and the shorter "if..then.." work without having to left factorize the parser cases, and there is no need to lookahead. They have no equivalent to the exception "**Stream.Error**": when all cases are tested, and have failed, the parsers return the value "**None**". Their drawback is that, when a parsing error happens, it is not easily possible to know the location of the error in the input, since the initial stream has not been modified: the system would indicate a failure at the first character of the first line: this is a general drawback of backtracking parsers. See the solutions found to this problem in the page about purely functional parsers.

## Chapter 7

# Functional parsers

Purely functional parsers are an alternative of stream parsers where the used stream type is a lazy non-destructive type : these streams are lazy values, like in classical stream parsers, but the values are not removed as long as the parsing advances. To make them work, the parsers of purely functional streams return, not the simple values, but a value of type `option` : `"Node"` meaning "no match" (the equivalent of the exception `"Parse.Failure"` of normal streams) and `"Some (r, s)"` meaning "the result is r and the remaining stream is s".

### 7.1 Syntax

The syntax of purely functional parsers, when loading `"pa_fstream.cmo"`, is the following:

```
expression ::= fparser
              | match-with-fparser
  fparser   ::= "fparser" pos-opt "[" parser-cases "]"
              | "fparser" pos-opt parser-case
match-with-fparser ::= "match" expression with "fparser" pos-opt "[" parser-cases "]"
                    | "match" expression with "fparser" pos-opt parser-case
parser-cases ::= parser-cases parser-case
              | <nothing>
parser-case  ::= "[:" stream-pattern ":" pos-opt "->" expression
              | "[:" ":" pos-opt "->" expression
stream-pattern ::= stream-patt-comp
                | stream-patt-comp ";" stream-pattern
stream-patt-comp ::= "'" pattern
                  | pattern "=" expression
                  | pattern
pos-opt          ::= pattern
                  | <nothing>
```

Notice that, in difference with classical parsers, there is no difference, in a stream pattern, between the first stream pattern component and the other ones. In particular, there is not this syntax "question mark" and expression optionnally ending those components. Moreover, the "lookahead" case is not necessary, we see further why. On the contrary, the syntaxes "pattern when" and "let..in" inside stream patterns we see in classical parsers are just not implemented.

## 7.2 Streams

The functional parsers are functions taking as parameters functional streams, which are values of type `"Fstream.t a"` for some type `"a"`. It is possible to build functional streams using the functions defined in the module `"Fstream"`:

### 7.2.1 Fstream.from

`"Fstream.from f"` returns a stream built from the function `"f"`. To create a new stream element, the function `"f"` is called with the current stream count, starting with zero. The user function `"f"` must return either `"Some <value>"` for a value or `"None"` to specify the end of the stream.

### 7.2.2 Fstream.of\_list

Return a stream built from the list in the same order.

### 7.2.3 Fstream.of\_string

Return a stream of the characters of the string parameter.

### 7.2.4 Fstream.of\_channel

Return a stream of the characters read from the input channel parameter.

## 7.3 Semantics of parsers

### 7.3.1 Fparser

The purely functional parsers act like classical parsers, with a recursive descent algorithm, except that:

- If the first stream pattern component matches the beginning of the stream, there is no error if the following stream patterns components do not match: the control simply passes to the next parser case with the initial stream.
- If the semantic actions are of type `"t"`, the result of the parser is of type `"option (t * Fstream.t)"`, not just `"t"` like in classical parsers. If a stream pattern matches, the semantic action is evaluated, giving some result `"e"` and the result of the parser is `"Some (e, strm)"` where `"strm"` is the remaining stream.
- If no parser case matches, the result of the parser is `"None"`.

### 7.3.2 Error position

A difficulty, with purely functional parsers, is how to find the position of the syntax error, when the input is wrong. Since the system tries all parsers cases before returning `"None"`, and that the initial stream is not affected, it is not possible to directly find where the error happened. This is a problem for parsing using backtracking (here, it is limited backtracking, but the problem is the same).

The solution is to use the function `"Fstream.count_unfrozen"` applied to the initial stream. Like its name says, it returns the number of unfrozen elements of the stream, which is exactly the longest match found. If the input is a stream of characters, the return of this function is exactly the position in number of characters from the beginning of the stream.

However, it is not possible to know directly which rule failed and therefore it is not possible, like in classical parsers, to specify and get clear error messages. Future versions of purely functional parsers may propose solutions to resolve this problem.

Notice that, if using that method, it is not possible to reuse the same stream to call another parser, and hope to get the right position of the error, if another error happens, since it may test less terminals than the first parser. Use a fresh stream in this case, if possible.



# Chapter 8

## Stream lexers

The file `pa_lex.cmo` is a Camlp5 syntax extension kit for parsers for streams of the type `'char'`. This syntax is shorter in length and more readable than its equivalent version written with classical stream parsers. Like classical parsers, they use recursive descendant parsing. They are also pure syntax sugar, and each lexer written with this syntax can be written using normal parsers syntax.

### 8.1 Introduction

Classical parsers in OCaml apply to streams of any type of values. But, for the specific type `"char"`, it has been possible to shorten codes in several ways, in particular by using strings to group several characters together, and by hiding the management of a `"lexing buffer"`, a data structure recording the matched characters.

In order to motivate our work, here is a small example. Let us look at the following parser. It parses an identifier, composed of letters, digits, underscores, quotes and utf-8 characters, and record the result in a buffer. In classical parsers syntax, this could be written like this:

```
value rec ident buf =
  parser
  [ [: ('A'..'Z' | 'a'..'z' | '0'..'9' | '_' | '\'' | '\128'..\255'
        as c); buf = ident (B.add c buf) ! :) -> buf
  | [: :] -> buf ]
;
```

With the new syntax, it is possible to write it as:

```
value rec ident = lexer [ "A..Za..z0..9_'\128..\255" ident! | ];
```

The two codes are strictly equivalent, but the lexer version is easier to write and understand, and is much shorter.

### 8.2 Syntax

When loading the syntax extension `pa_lex.cmo`, the OCaml syntax is extended as follows:

```
expression ::= lexer
lexer ::= "lexer" "[" rules "]"
rules ::= rules rule
```

```

        | <nothing>
    rule ::= symbols [ ">" action ]
symbols ::= symbols symbol err
        | <nothing>
symbol  ::= "_" no-record-opt
        | STRING no-record-opt
        | simple-expression
        | "?=" "[" lookaheads "]"
        | "[" rules "]"
no-record-opt ::= "/"
        | <nothing>
simple-expression ::= LIDENT
        | CHAR
        | "(" <expression> ")"
lookaheads ::= lookaheads "|" lookahead-sequence
        | lookahead-sequence
lookahead-sequence ::= lookahead-symbols
        | STRING
lookahead-symbols ::= lookahead-symbols lookahead-symbol
        | lookahead-symbol
lookahead-symbol ::= CHAR
        | "_"
err ::= "?" simple-expression
        | "!"
        | <nothing>
action ::= expression

```

The identifiers "STRING", "CHAR" and "LIDENT" above represent the OCaml tokens corresponding to string, character and lowercase identifier (identifier starting with a lowercase character).

Moreover, together with that syntax extension, another extension is added the entry **expression**, typically for the semantics actions of the "lexer" statement above, but not only. It is:

```

expression ::= "$" "add" STRING
        | "$" "buf"
        | "$" "empty"
        | "$" "pos"

```

Remark: the identifiers "add", "buf", "empty" and "pos" do not become keywords (they are not reserved words) but just identifiers. On the contrary, the identifier **lexer**, which introduces the syntax, is a new keyword and cannot be used as variable identifier any more.

## 8.3 Semantics

A lexer defined in the syntax above is a shortcut version of a parser applied to the specific case of streams of characters. It could be written with a normal parser. The proposed syntax is much shorter, easier to use and to understand, and silently takes care of the lexing buffer for the programmer. The lexing buffers are data structures, which are passed as parameters to called lexers and returned by them.

Our lexers are of the type:

```
B.t -> Stream.t char -> u
```

where "u" is a type which depends on what the lexer returns. If there is no semantic action (since it is optional), this type is automatically "B.t".

"B" is a module which must be defined by the user. It has to contain the lexing buffer type "t" and some variables and functions:

- **empty**: the empty lexing buffer
- **add**: the way to add a character to a lexing buffer
- **get**: the way to get a string from the lexing buffer

A possible implementation, using "list char" as lexing buffer type ("B.t"), recording the characters at top of the list (therefore creating a list in reverse order) could be:

```
(* tool function, converting a reversed list of char into a string *)
value rev_implode cl =
  let s = String.create (List.length cl) in
  loop (String.length s - 1) cl where rec loop i =
    fun
      [ c :: cl ] -> do { s.[i] := c; loop (i - 1) cl }
    | [] -> s ]
;

(* the lexing buffer module *)
module B =
  struct
    type t = list char;
    value empty = [];
    value add c l = [c :: l];
    value get = rev_implode;
  end
;
```

A lexer is a function with two parameters: the first one is the lexing buffer itself, and the second one the stream. When called, it tries to match the stream against its first rule. If it fails, it tries its second rule, and so on, up to its last rule. If the last rule fails, the lexer fails by raising the exception "Stream.Failure". All of this is the usual behaviour of stream parsers.

In a rule, when a character is matched, it is inserted into the lexing buffer, except if the "no record" feature is used (see further).

Rules which have no semantic action return the lexing buffer itself.

### 8.3.1 Symbols

The different kinds of symbols in a rule are:

- The token "underscore", which represents any character. Fails only if the stream is empty.
- A string which represents a matching of any character in the string. Notice that it is a choice between all the characters, not the sequence of these characters. To indicate a sequence, you have to use several symbols, the ones behind the others. Character ranges can be inserted using the characters "...". For example, to specify a match of any letter or digit, you can write "A..Za..z0..9". Conversely, the

beginning of a comment in the OCaml language has to be written: `"(" *`, not `"(*"` which would mean "the character left parenthesis or the character star".

- An expression corresponding to a call to another lexer, which takes the buffer as first parameter and has to return another lexing buffer with all characters found in the stream catenated to the lexing buffer.
- The sequence `"?="` introducing lookahead characters. See the section 3.3.
- A rule, recursively, between brackets, inlining a lexer.

In the first two cases (namely, underscore and string), the symbol can be optionally followed by the character `"slash"` specifying that the found symbol must not be added into the lexing buffer. By default, it is. Useful, for example, when writing a lexer parsing strings, when the initial double quote and final double quote have not to be part of the string itself.

Moreover, a symbol can be followed by an optional error indicator, which can be:

- The character `?` (question mark) followed by a string expression, telling that, if there is a syntax error at this point (i.e. the symbol is not matched although the beginning of the rule was), the exception `Stream.Error` is raised with that string as parameter. Without this indicator, it is raised with the empty string. This is the same behaviour than with classical stream parsers.
- The character `!` (exclamation mark), which is just an indicator to let the syntax expander optimize the code. If the programmer is sure that the symbol never fails (i.e. never raises `Stream.Failure`), in particular if this symbol recognizes the empty rule, he can add this exclamation mark. If it is used correctly (the compiler cannot check it), the behaviour is identical as without the `!`, except that the code is shorter and faster, and can sometimes be tail recursive. If the indication is not correct, the behaviour of the lexer is undefined. This feature exists also in classical stream parsers (it is a new feature added in 2007).

### 8.3.2 Specific expressions

When loading this syntax extension, the entry `<expression>`, at level labelled "simple" of the OCaml language is extended with the following rules:

- `$add` followed by a string, specifying that the programmer wants to add all characters of the string in the lexing buffer. It returns the new lexing buffer. It corresponds to an iteration of calls to `B.add` with all characters of the string with the current lexing buffer as initial parameter.
- `$buf` which returns the lexing buffer converted into string.
- `$empty` which returns an empty lexing buffer.
- `$pos` which return the current position in the stream in number of characters (starting at zero).

### 8.3.3 Lookahead

Lookahead is useful in some cases, when factorization of rules is impossible. To understand how it is useful, a first remark must be done, about the usual behaviour of Camlp5 stream parsers.

Stream parsers (including these lexers) use a limited parsing algorithm, in a way that when the first symbol of a rule is matched, all the following symbols of the same rule must apply, otherwise it is a syntax error. There is no backtrack. In most of the cases, left factorization of rules resolve conflicting problems. For

example, in parsers of tokens (which is not our case here, since we parse only characters), when one writes a parser to recognize both the typical grammar rules "if..then..else" and the shorter "if..then..", the solution is to write a rule starting with "if..then.." followed with a call to a parser recognizing "else.." or nothing.

Sometimes, however, this left factorization is not possible. A lookahead of the stream to check the presence of some elements (these element being characters, if we are using this "lexer" syntax) might be necessary to decide whether or not it is a good idea to start the rule. This lookahead feature may unfreeze several characters from the input stream but without removing them.

Syntactically, a lookahead starts with `?=` and is followed by one or several lookahead sequences separated by the vertical bar `|`, the whole list being enclosed by braces.

If there are several lookaheads, they must all be of the same size (contain the same number of characters).

If the lookahead sequence is just a string, it corresponds to all characters of this string in the order (which is different for strings outside lookahead sequences, representing a choice of all characters).

Examples of lookaheads:

```
?= [ _ ' ' | '\ ' _ ]  
?= [ "<<" | "<:" ]
```

The first line above matches a stream whose second character is a quote or a stream whose first character is a backslash (real example in the lexer of OCaml, in the library of Camlp5, named "plexer.ml"). The second line matches a stream starting with the two characters `<` and `<` or starting with the two characters `<` and `:` (this is another example in the same file).

### 8.3.4 Semantic actions of rules

By default, the result of a "lexer" is the current lexing buffer, which is of type `B.t`. But it is possible to return other values, by adding `"->"` at end of rules followed by the expression you want to return, like in usual pattern matching in OCaml.

An interesting result, for example, could be the string corresponding to the characters of the lexing buffer. This can be obtained by returning the value `"$buf"` (see section 3.2)

### 8.3.5 A complete example

A complete example can be seen in the sources of Camlp5, file `"lib/plexer.ml"`. This is the lexer of OCaml, either "normal" or "revised" syntax.

### 8.3.6 Compiling

To compile a file containing lexers, just load `pa_lex.cmo` using one of the following methods:

- Either by adding `pa_lex.cmo` among the camlp5 options. See the camlp5 manual page or documentation.
- Or by adding `#load "pa_lex.cmo";` anywhere in the file, before the usages of this "lexer" syntax.

### 8.3.7 How to display the generated code

You can see the generated code, for a file "bar.ml" containing lexers, by typing in a command line:

```
camlp5r pa_lex.cmo pr_r.cmo bar.ml
```

To see the equivalent code with stream parsers, use:

```
camlp5r pa_lex.cmo pr_r.cmo pr_rp.cmo bar.ml
```

## Chapter 9

# Extensible grammars

This chapter describes the whole syntax and semantics of the extensible grammars of `camlp5`.

The extensible grammars are the most advanced parsing tool of `camlp5`. They apply to streams of characters using a lexer which has to be previously defined by the programmer. In `camlp5`, the syntax of the `ocaml` language is defined with extensible grammars, which makes `camlp5` a bootstrapped system (it compiles its own features by itself).

### 9.1 Getting started

The extensible grammars are a system to build *grammar entries* which can be extended dynamically. A grammar entry is an abstract value internally containing a stream parser. The type of a grammar entry is `"Grammar.Entry.e t"` where `"t"` is the type of the values returned by the grammar entry.

To start with extensible grammars, it is necessary to build a *grammar*, a value of type `"Grammar.g"`, using the function `"Grammar.gcreate"`:

```
value g = Grammar.gcreate lexer;
```

where `"lexer"` is a lexer previously defined. See the section explaining the interface with lexers. In a first time, it is possible to use a lexer of the module `"Plexer"` provided by `camlp5`:

```
value g = Grammar.gcreate (Plexer.gmake ());
```

Each grammar entry is associated with a grammar. Only grammar entries of the same grammar can call each other. To create a grammar entry, one has to use the function `"Grammar.Entry.create"` which takes the grammar as first parameter and a name as second parameter. This name is used in case of syntax errors. For example:

```
value exp = Grammar.Entry.create g "expression";
```

To apply a grammar entry, the function `"Grammar.Entry.parse"` can be used. Its first parameter is the grammar entry, the second one a stream of characters:

```
Grammar.Entry.parse exp (Stream.of_string "hello");
```

But if you experiment this, since the entry was just created without any rules, you receive an error message:

```
Stream.Error "entry [expression] is empty"
```

To add grammar rules to the grammar entry, it is necessary to *extend* it, using a specific syntactic statement: `"EXTEND"`.

## 9.2 Syntax of the EXTEND statement

The "EXTEND" statement is added in the expressions of the ocaml language when the syntax extension kit "pa\_extend.cmo" is loaded. Its syntax is:

```
expression ::= extend
  extend ::= "EXTEND" extend-body "END"
extend-body ::= global-opt entries
global-opt ::= "GLOBAL" ":" entry-names ";"
              | <nothing>
entry-names ::= entry-name entry-names
              | entry-name
  entry ::= entry-name ":" position-opt "[" levels "]"
position-opt ::= "FIRST"
              | "LAST"
              | "BEFORE" label
              | "AFTER" label
              | "LEVEL" label
              | <nothing>
levels ::= level "|" levels
        | level
level ::= label-opt assoc-opt "[" rules "]"
label-opt ::= label
           | <nothing>
assoc-opt ::= "LEFTA"
            | "RIGHTA"
            | "NONA"
            | <nothing>
rules ::= rule "|" rules
        | rule
rule ::= psymbols-opt "->" expression
      | psymbols-opt
psymbols-opt ::= psymbols
              | <nothing>
psymbols ::= psymbol ";" psymbols
           | psymbol
psymbol ::= symbol
          | pattern "=" symbol
symbol ::= keyword
        | token
        | token string
        | entry-name
        | entry-name "LEVEL" label
        | "SELF"
        | "NEXT"
        | "LIST0" symbol
        | "LIST0" symbol "SEP" symbol
        | "LIST1" symbol
        | "LIST1" symbol "SEP" symbol
        | "OPT" symbol
        | "[" rules "]"
        | "(" symbol ")"
```

```

keyword ::= string
token  ::= uident
label  ::= string
entry-name ::= qualid
qualid ::= qualid "." qualid
        | uident
        | lident
uident ::= 'A'-'Z' ident
lident ::= ('a'-'z' | '_' | utf8-char) ident
ident  ::= ident-char*
ident-char ::= ('a'-'a' | 'A'-'Z' | '0'-'9' | '_' | '\'' | utf8-char)
utf8-char ::= '\128'-''\255'

```

Other statements, "GEXTEND", "DELETE\_RULE", "GDELETE\_RULE" are also defined by the same syntax extension kit. See further.

In the description above, only "EXTEND" and "END" are new keywords (reserved words which cannot be used in variables, constructors or module names). The other strings (e.g. "GLOBAL", "LEVEL", "LISTO", "LEFTA", etc.) are not reserved.

## 9.3 Semantics of the EXTEND statement

The EXTEND statement starts with the "EXTEND" keyword and ends with the "END" keyword.

### 9.3.1 GLOBAL indicator

After the first keyword, it is possible to see the identifier "GLOBAL" followed by a colon, a list of entries names and a semicolon. It says that these entries correspond to visible (previously defined) entry variables, in the context of the EXTEND statement, the other ones being locally and silently defined inside.

- If an entry, which is extended in the EXTEND statement, is in the GLOBAL list, but is not defined in the context of the EXTEND statement, the ocaml compiler will fail with the error "unbound value".
- If there is no GLOBAL indicator, and an entry, which is extended in the EXTEND statement, is not defined in the context of the EXTEND statement, the ocaml compiler will also fail with the error "unbound value".

Example:

```

value exp = Grammar.Entry.create g "exp";
EXTEND
  GLOBAL: exp;
  exp: [ [ x = foo; y = bar ] ];
  foo: [ [ "foo" ] ];
  bar: [ [ "bar" ] ];
END;

```

The entry "exp" is an existing variable (defined by value exp = ...). On the other hand, the entries "foo" and "bar" have not been defined. Because of the GLOBAL indicator, the system defines them locally.

Without the GLOBAL indicator, the three entries would have been considered as global variables, therefore the ocaml compiler would say "unbound variable" under the first undefined entry, "foo".

### 9.3.2 Entries list

Then the list of entries extensions follow. An entry extension starts with the entry name followed by a colon. An entry may have several levels corresponding to several stream parsers which call the ones the others (see further).

#### Optional position

After the colon, it is possible to specify a where to insert the defined levels:

- The identifier "FIRST" (resp. "LAST") indicates that the level must be inserted before (resp. after) all possibly existing levels of the entry. They become their first (resp. last) levels.
- The identifier "BEFORE" (resp. "AFTER") followed by a level label (a string) indicates that the levels must be inserted before (resp. after) that level, if it exists. If it does not exist, the extend statement fails at run time.
- The identifier "LEVEL" followed by a level label indicates that the first level defined in the extend statement must be inserted at the given level, extending and modifying it. The other levels defined in the statement are inserted after this level, and before the possible levels following this level. If there is no level with this label, the extend statement fails at run time.
- By default, if the entry has no level, the levels defined in the statement are inserted in the entry. Otherwise the first defined level is inserted at the first level of the entry, extending or modifying it. The other levels are inserted afterwards (before the possible second level which may previously exist in the entry).

#### Levels

After the optional "position", the *level* list follow. The levels are separated by vertical bars, the whole list being between brackets.

A level starts with an optional label, which corresponds to its name. This label is useful to specify this level in case of future extension, using the *position* (see previous section).

The level continues with an optional associativity indicator, which can be:

- LEFTA for left associativity (default),
- RIGHTA for right associativity,
- NONA for no associativity.

#### Rules

At last, the grammar *rule* list appear. The rules are separated by vertical bars, the whole list being brackets.

A rule looks like a match case in the "match" statement or a parser case in the "parser" statement: a list of psymbols (see next paragraph) separated by semicolons, followed by a right arrow and an expression, the semantic action. Actually, the right arrow and expression are optional: in this case, it is equivalent to an expression which would be the unit "()" constructor.

A psymbol is either a pattern, followed with the equal sign and a symbol, or by a symbol alone. It corresponds to a test of this symbol, whose value is bound to the pattern if any.

## Symbols

A symbol is either:

- a keyword (a string): the input must match this keyword,
- a token name (an identifier starting with an uppercase character), optionally followed by a string: the input must match this token (any value if no string, or that string if a string follows the token name), the list of the available tokens depending on the associated lexer (the list of tokens available with "Plexer.gmake ()" is: LIDENT, UIDENT, TILDEIDENT, TILDEIDENTCOLON, QUESTION-IDENT, INT, INT\_1, INT\_L, INT\_n, FLOAT, CHAR, STRING, QUOTATION, ANTIQUOT and EOF; other lexers may propose other lists of tokens),
- an entry name, which correspond to a call to this entry,
- an entry name followed by the identifier "LEVEL" and a level label, which correspond to the call to this entry at that level,
- the identifier "SELF" which is a recursive call to the present entry, according to the associativity (i.e. it may be a call at the current level, to the next level, or to the top level of the entry): "SELF" is equivalent to the name of the entry itself,
- the identifier "NEXT", which is a call to the next level of the current entry,
- a left brace, followed by a list of rules separated by vertical bars, and a right brace: equivalent to a call to an entry, with these rules, inlined,
- a meta symbol (see further),
- a symbol between parentheses.

The syntactic analysis follow the list of symbols. If it fails, depending on the first items of the rule (see the section about the kind of grammars recognized):

- the parsing may fail by raising the exception "**Stream.Error**"
- the parsing may continue with the next rule.

## Meta symbols

Extra symbols exist, allowing to manipulate lists or optional symbols. They are:

- LIST0 followed by a symbol: this is a list of this symbol, possibly empty,
- LIST0 followed by a symbol, SEP and another symbol: this is a list, possibly empty, of the first symbol separated by the second one,
- LIST1 followed by a symbol: this is a list of this symbol, with at least one element,
- LIST0 followed by a symbol, SEP and another symbol: this is a list, with at least one element, of the first symbol separated by the second one,
- OPT followed by a symbol: equivalent to "this symbol or nothing".

### 9.3.3 Rules insertion

Remember that "EXTEND" is a statement, not a declaration: the rules are added in the entries at run time. Each rule is internally inserted in a tree, allowing the left factorization of the rule. For example, with this list of rules (borrowed from the camlp5 sources):

```
"method"; "private"; "virtual"; l = label; ":"; t = poly_type
"method"; "virtual"; "private"; l = label; ":"; t = poly_type
"method"; "virtual"; l = label; ":"; t = poly_type
"method"; "private"; l = label; ":"; t = poly_type; "="; e = expr
"method"; "private"; l = label; sb = fun_binding
"method"; l = label; ":"; t = poly_type; "="; e = expr
"method"; l = label; sb = fun_binding
```

the rules are inserted in a tree and the result looks like:

```
"method"
|-- "private"
|   |-- "virtual"
|   |   |-- label
|   |   |-- ":"
|   |   |-- poly_type
|   |-- label
|   |   |-- ":"
|   |   |-- poly_type
|   |   |-- "!="
|   |   |-- expr
|   |-- fun_binding
|-- "virtual"
|   |-- "private"
|   |   |-- label
|   |   |-- ":"
|   |   |-- poly_type
|   |-- label
|   |   |-- ":"
|   |   |-- poly_type
|-- label
|   |-- ":"
|   |   |-- poly_type
|   |   |-- "="
|   |   |-- expr
|   |-- fun_binding
```

This tree is built as long as rules are inserted. When used, by applying the function "Grammar.Entry.parse" to the current entry, the input is matched with that tree, starting from the tree root, descending on it as long as the parsing advances.

There is a different tree by entry level.

### 9.3.4 Semantic action

The semantic action, i.e. the expression following the right arrow in rules, contain in its environment:

- the variables bound by the patterns of the symbols found in the rules,

- the specific variable "loc" which contain the location of the whole rule in the source.

The location is an abstract type defined in the module "Stdpp" of camlp5.

It is possible to change the name of this variable by using the option "-loc" of camlp5. For example, compiling a file like this:

```
camlp5r -loc foobar file.ml
```

the variable name, for the location will be "foobar" instead of "loc".

## 9.4 The DELETE\_RULE statement

The "DELETE\_RULE" statement is also added in the expressions of the ocaml language when the syntax extension kit "pa\_extend.cmo" is loaded. Its syntax is:

```
expression ::= delete-rule
delete-rule ::= "DELETE_RULE" delete-rule-body "END"
delete-rule-body ::= entry-name ":" symbols
                  symbols ::= symbol symbols
                           | symbol
```

See the syntax of the EXTEND statement for the meaning of the syntax entries not defined above.

The entry is scanned for a rule matching the giving symbol list. When found, the rule is removed. If no rule is found, the exception "Not\_found" is raised.

## 9.5 Extensions FOLD0 and FOLD1

...To be written...

## 9.6 Extensions SLIST0, SLIST1 and SOPT

...To be written...

## 9.7 Grammar machinery

We explain here the detail of the mechanism of the parsing of an entry.

### 9.7.1 Start and Continue

At each entry level, the rules are separated into two trees:

- The tree of the rules *not* starting with the current entry name nor by "SELF".
- The tree of the rules starting with the current entry name or by the identifier "SELF", this symbol not being included in the tree.

They determine two functions:

- The function named "start", analyzing the first tree.

- The function named "continue", taking, as parameter, a value previously parsed, and analyzing the second tree.

A call to an entry, using "Grammar.Entry.parse" correspond to a call to the "start" function of the first level of the entry.

The "start" function tries its associated tree. If it works, it calls the "continue" function of the same level, giving the result of "start" as parameter. If this "continue" function fails, this parameter is simply returned. If the "start" function fails, the "start" function of the next level is tested. If there is no more levels, the parsing fails.

The "continue" function first tries the "continue" function of the next level. If it fails, or if it is the last level, it tries its associated tree, then calls itself again, giving the result as parameter. If its associated tree fails, it returns its extra parameter.

## 9.7.2 Associativity

While testing the tree, there is a special case for rules ending with SELF or with the current entry name. For this last symbol, there is a call to the "start" function: of the current level if the level is right associative, or of the next level otherwise.

There is no behaviour difference between left and non associative, because, in case of syntax error, the system attempts, anyway, to recover the error by applying the "continue" function of the previous symbol (if this symbol is a call to an entry).

When a SELF or the current entry name is encountered in the middle of the rule (i.e. if it is not the last symbol), there is a call to the "start" function of the first level of the current entry.

Example. Let us consider the following grammar:

```
EXTEND
  expr:
    [ "minus" LEFTA
      [ x = SELF; "-"; y = SELF -> x -. y ]
    | "power" RIGHTA
      [ x = SELF; "**"; y = SELF -> x ** y ]
    | "simple"
      [ "("; x = SELF; ")" -> x
      | x = INT -> float_of_int x ] ]
;
```

The left "SELF"s of the two levels "minus" and "power" correspond to a call to the next level. In the level "minus", the right "SELF" also, and the left associativity is treated by the fact that the "continue" function is called (starting with the keyword "-" since the left "SELF" is not part of the tree). On the other hand, for the level "power", the right "SELF" corresponds to a call to the current level, i.e. the level "power" again. At end, the "SELF" between parentheses of the level "simple" correspond to a call to the first level, namely "minus" in this grammar.

## 9.7.3 Errors and recovery

Like for stream parsers, two exceptions may happen: "Stream.Failure" or "Stream.Error". The first one indicates that the parsing just could not start. The second one indicates that the parsing started but failed further.

In stream parsers, when the first symbol of a rule has been accepted, all the symbols of the same rule must be accepted, otherwise the exception "Stream.Error" is raised.

Here, in extensible grammars, unlike stream parsers, before the "Stream.Error" exception, the system attempts to recover the error by the following trick: if the previous symbol of the rule was a call to another entry, the system calls the "continue" function of that entry, which may resolve the problem.

In extensible grammars, the exceptions are encapsulated with the exception "Stdpp.Exc\_located" giving the location of the error together with the exception itself.

#### 9.7.4 Tokens starting rules

Another improvement (than the error recovery) is the fact that, when a rule starts with several tokens and/or keywords, all these tokens and keywords are tested in one time, and the possible "Stream.Error" may happen, only from the symbol following them on, if any.

#### 9.7.5 Kind of grammar

The kind of grammar is predictive parsing grammar, i.e. recursive descent parsing without backtrack. But with some nuances, due to the improvements (error recovery and token starting rules) indicated in the previous sections.

### 9.8 The Grammar module

The Grammar module contains all what is necessary to manipulate grammars and entries. It contains:

#### 9.8.1 Main types and values

```
type g = 'abstract;
```

The type of grammars, holding entries.

```
value gcreate : Token.glexer (string * string) -> g;
```

Create a new grammar, without keywords, using the lexer given as parameter.

```
value tokens : g -> string -> list (string * int);
```

Given a grammar and a token pattern constructor, returns the list of the corresponding values currently used in all entries of this grammar. The integer is the number of times this pattern value is used.

Examples:

- The call [Grammar.tokens g ""] returns the keywords list.
- The call [Grammar.tokens g "IDENT"] returns the list of all usages of the pattern "IDENT" in the EXTEND statements.

```
value gllexer : g -> Token.glexer token;
```

Return the lexer used by the grammar

```
type parsable = 'abstract;
```

```
value parsable : g -> Stream.t char -> parsable;
```

Type and value allowing to keep the same token stream between several calls of entries of the same grammar, to prevent loss of tokens. To be used with Entry.parse\_parsable below

```

module Entry =
  sig
    type e 'a = 'x;
    value create : g -> string -> e 'a;
    value parse : e 'a -> Stream.t char -> 'a;
    value parse_token : e 'a -> Stream.t token -> 'a;
    value parse_parsable : e 'a -> parsable -> 'a;
    value name : e 'a -> string;
    value of_parser : g -> string -> (Stream.t token -> 'a) -> e 'a;
    value print : e 'a -> unit;
    value find : e 'a -> string -> e Obj.t;
    external obj : e 'a -> Gramext.g_entry token = "%identity";
  end;

```

Module to handle entries.

- /Entry.e/ is the type for entries returning values of type [*a*].
- /Entry.create g n/ creates a new entry named [*n*] in the grammar [*g*].
- /Entry.parse e/ returns the stream parser of the entry [*e*].
- /Entry.parse\_token e/ returns the token parser of the entry [*e*].
- /Entry.parse\_parsable e/ returns the parsable parser of the entry [*e*].
- /Entry.name e/ returns the name of the entry [*e*].
- /Entry.of\_parser g n p/ makes an entry from a token stream parser.
- /Entry.print e/ displays the entry [*e*] using [Format].
- /Entry.find e s/ finds the entry named [*s*] in [*e*]'s rules.
- /Entry.obj e/ converts an entry into a [Gramext.g\_entry] allowing to see what it holds ([Gramext] is visible, but not documented).

```

value of_entry : Entry.e 'a -> g;

```

Return the grammar associated with an entry.

## 9.8.2 Printing grammar entries

The function "Grammar.Entry.print" displays the current contents of an entry. Interesting for debugging, to look at the result of a syntax extension, to see the names of the levels.

The display does not include the patterns nor the semantic actions, whose sources are not recorded in the grammar entries data.

Moreover, the local entries (see the section about the GLOBAL indicator) are indicated with a star ("\*") to inform that they are not directly accessible.

## 9.8.3 Clearing grammars and entries

```

module Unsafe :
  sig
    value gram_reinit : g -> Token.glexer token -> unit;
    value clear_entry : Entry.e 'a -> unit;
  end;

```

Module for clearing grammars and entries. To be manipulated with care, because: 1) reinitializing a grammar destroys all tokens and there may have problems with the associated lexer if there are keywords; 2) clearing an entry does not destroy the tokens used only by itself.

- /Unsafe.reinit\_gram g lex/ removes the tokens of the grammar and sets [lex] as a new lexer for [g]. Warning: the lexer itself is not reinitialized.
- /Unsafe.clear\_entry e/ removes all rules of the entry [e].

## 9.8.4 Functorial interface

Alternative for grammars use. Grammars are no more Ocaml values: there is no type for them. Modules generated preserve the rule "an entry cannot call an entry of another grammar" by normal OCaml typing.

```
module type GLexerType =
  sig
    type te = 'x;
    value lexer : Token.glexer te;
  end;
```

The input signature for the functor [Grammar.GMake]: [te] is the type of the tokens.

```
module type S =
  sig
    type te = 'x;
    type parsable = 'x;
    value parsable : Stream.t char -> parsable;
    value tokens : string -> list (string * int);
    value glexer : Token.glexer te;
    module Entry :
      sig
        type e 'a = 'y;
        value create : string -> e 'a;
        value parse : e 'a -> parsable -> 'a;
        value parse_token : e 'a -> Stream.t te -> 'a;
        value name : e 'a -> string;
        value of_parser : string -> (Stream.t te -> 'a) -> e 'a;
        value print : e 'a -> unit;
        external obj : e 'a -> Gramext.g_entry te = "%identity";
      end;
    module Unsafe :
      sig
        value gram_reinit : Token.glexer te -> unit;
        value clear_entry : Entry.e 'a -> unit;
      end;
    value extend :
      Entry.e 'a -> option Gramext.position ->
        list
          (option string * option Gramext.g_assoc *
            list (list (Gramext.g_symbol te) * Gramext.g_action)) ->
            unit;
    value delete_rule : Entry.e 'a -> list (Gramext.g_symbol te) -> unit;
  end;
```

Signature type of the functor [Grammar.GMake]. The types and functions are almost the same than in generic interface, but:

- Grammars are not values. Functions holding a grammar as parameter do not have this parameter yet.
- The type [parsable] is used in function [parse] instead of the char stream, avoiding the possible loss of tokens.
- The type of tokens (expressions and patterns) can be any type (instead of (string \* string)); the module parameter must specify a way to show them as (string \* string).

```
module GMake (L : GLexerType) : S with type te = L.te;
```

### 9.8.5 Miscellaneous

```
value error_verbose : ref bool;
```

Flag for displaying more information in case of parsing error; default = [False].

```
value warning_verbose : ref bool;
```

Flag for displaying warnings while extension; default = [True].

```
value strict_parsing : ref bool;
```

Flag to apply strict parsing, without trying to recover errors; default = [False].

```
value print_entry : Format.formatter -> Gramext.g_entry 'te -> unit;
```

General printer for all kinds of entries (obj entries).

```
value iter_entry : (Gramext.g_entry 'te -> unit) -> Gramext.g_entry 'te -> unit;
```

[Grammar.iter\_entry f e] applies [f] to the entry [e] and transitively all entries called by [e]. The order in which the entries are passed to [f] is the order they appear in each entry. Each entry is passed only once. \*)

```
value fold_entry : (Gramext.g_entry 'te -> 'a -> 'a) -> Gramext.g_entry 'te -> 'a -> 'a;
```

[Grammar.fold\_entry f e init] computes [(f eN .. (f e2 (f e1 init)))], where [e1 .. eN] are [e] and transitively all entries called by [e]. The order in which the entries are passed to [f] is the order they appear in each entry. Each entry is passed only once. \*)

## 9.9 Interface with the lexer

To create a grammar, the function "Grammar.gcreate" must be called, with a lexer as parameter.

A simple solution, as possible lexer, is the predefined lexer built by "Plexer.gmake ()", lexer used for the ocaml grammar of camlp5. In this case, you can just put it as parameter of "Grammar.gcreate" and it is not necessary to read this section.

The section first introduces the notion of "token patterns" which are the way the tokens and keywords symbols in the EXTEND statement are represented. Then follow the description of the type of the parameter of "Grammar.gcreate".

### 9.9.1 Token patterns

A token pattern is a value of the type defined like this:

```
type pattern = (string * string);
```

This type represents values of the token and keywords symbols in the grammar rules.

For a token symbol in the grammar rules, the first string is the token constructor name (starting with an uppercase character), the second string indicates whether the match is "any" (the empty string) or some specific value of the token (an non-empty string).

For a keyword symbol, the first string is empty and the second string is the keyword itself.

For example, given this grammar rule:

```
"for"; i = LIDENT; "="; e1 = SELF; "to"; e2 = SELF
```

the different symbols and keywords are represented by the following couples of strings:

- the keyword "for" is represented by ("", "for"),
- the keyword "=" by ("", "="),
- the keyword "to" by ("", "to"),
- and the token symbol LIDENT by ("LIDENT", "").

The symbol UIDENT "Foo" in a rule would be represented by the token pattern:

```
("UIDENT", "Foo")
```

Notice that the symbol "SELF" is a specific symbol of the EXTEND syntax: it does not correspond to a token pattern and is represented differently. A token constructor name must not belong to the specific symbols: SELF, NEXT, LIST0, LIST1 and OPT.

### 9.9.2 The gllexer record

The type of the parameter of the function "Grammar.gcreate" is "gllexer", defined in the module "Token". It is a record type with the following fields:

```
tok_func
```

It is the lexer itself. Its type is:

```
Stream.t char -> (Stream.t (string * string) * location_function);
```

The lexer takes a character stream as parameter and must answer a couple of: a token stream, the tokens being represented by a couple of strings, and a location function.

The location function is a function taking, as parameter, a integer corresponding to a token number in the stream (starting from zero), and returning the location of this token in the source. It is important to get the good locations in the semantic actions of the grammar rules.

Notice that, despite the lexer takes a character stream as parameter, it is not mandatory to use the stream parsers technology to write the lexer. What is important is that it does the job.

`tok_using`

It is a function of type:

```
pattern -> unit
```

The parameter of this function is the representation of a token symbol or a keyword symbol in grammar rules. See the section about token patterns.

This function is called for each token symbol and each keyword encountered in the grammar rules of the EXTEND statement. Its goal is to allow the lexer to check that the tokens and keywords do respect the lexer rules. It checks that the tokens exist and are not misspelled. It can be also used to enter the keywords in the lexer keyword tables.

Setting it as the function that does nothing is possible, but the check of correctness of tokens is not done.

In case of error, the function must raise the exception `"Token.Error"` with an error message as parameter.

`tok_removing`

It is a function of type:

```
pattern -> unit
```

It is called by the DELETE\_RULE statement for each token symbol and each keyword that an occurrence of them is no more used. This can be interesting for keywords, if the lexer record the number of occurrences of the keywords: when the number of occurrences falls to zero, the keyword can be removed from the lexer tables.

`tok_match`

It is a function of type:

```
pattern -> ((string * string) -> unit)
```

The function tells how a token of the input stream is matched against a token pattern. Both are represented by a couple of strings.

This function takes a token pattern as parameter and return a function matching a token, returning the matched string or raising the exception `"Stream.Failure"` if the token does not match.

Notice that, for efficiency, it is necessary to write this function as a match of token patterns returning, for each case, the function which matches the token, *not* a function matching the token pattern and the token together and returning a string for each case.

An acceptable function is provided in the module `"Token"` and is named `"default_match"`. Its code looks like this:

```
value default_match =
  fun
  [ (p_con, "") ->
    fun (con, prm) -> if con = p_con then prm else raise Stream.Failure
  | (p_con, p_prm) ->
    fun (con, prm) ->
      if con = p_con && prm = p_prm then prm else raise Stream.Failure ]
;
```

`tok_text`

It is a function of type:

```
pattern -> string
```

Destinated to error messages, it takes a token pattern as parameter and return the string giving its name.

It is possible to use the predefined function `"lexer_text"` of the Token module. This function just returns the name of the token pattern constructor and its parameter if any.

For example, with this default function, the token symbol IDENT would be written as IDENT in error message (e.g. "IDENT expected"). The `"text"` function may decide to print it differently, e.g., as "identifier".

`tok_comm`

It is a mutable field of type:

```
option (list location)
```

It asks the lexer (the lexer function should do it) to record the locations of the comments in the program. Setting this field to "None" indicates that the lexer must not record them. Setting it to "Some []" indicated that the lexer must put the comments location list in the field, which is mutable.

### 9.9.3 Minimalist version

If a lexer have been written, named `"lexer"`, here is the minimalist version of the value suitable as parameter to `"Grammar.gcreate"`:

```
{Token.tok_func = lexer;
 Token.tok_using _ = (); Token.tok_removing _ = ();
 Token.tok_match = Token.default_match;
 Token.tok_text = Token.lexer_text;
 Token.tok_comm = None}
```

## 9.10 Functorial interface

The normal interface for grammars described in the previous sections has two drawbacks:

- First, the type of tokens of the lexers must be `"(string * string)"`
- Second, since the entry type has no parameter to specify the grammar it is bound to, there is no static check that entries are compatible, i.e. belong to the same grammar. The check is done at run time.

The functorial interface resolve these two problems. The functor takes a module as parameter where the token type has to be defined, together with the lexer returning streams of tokens of this type. The resulting module define entries compatible the ones to the other, and this is controlled by the ocaml type checker.

The syntax extension must be done with the statement `GEXTEND`, instead of `EXTEND`, and deletion by `GDELETE_RULE` instead of `DELETE_RULE`.

### 9.10.1 The gllexer type

In the section about the interface with the lexer, we presented the gllexer type as a record without type parameter. Actually, this type is defined as:

```
type gllexer 'te =
  { tok_func : lexer_func 'te;
    tok_using : pattern -> unit;
    tok_removing : pattern -> unit;
    tok_match : pattern -> 'te -> string;
    tok_text : pattern -> string;
    tok_comm : mutable option (list location) }
;
```

where the type parameter is the type of the token, which can be any type, different from "(string \* string)", providing the lexer function (tok\_func) returns a stream of this token type and the match function (tok\_match) indicates how to match values of this token type against the token patterns (which remain defined as "(string \* string)").

Here is an example of an user token type and the associated match function:

```
type mytoken = [ Ident of string | Int of int | Comma | Equal | Keyw of string ];

value mymatch =
  fun
  [ ("IDENT", "") -> fun [ Ident s -> s | _ -> raise Stream.Failure ]
  | ("INT", "") -> fun [ Int i -> string_of_int i | _ -> raise Stream.Failure ]
  | ("", ",") -> fun [ Comma -> "" | _ -> raise Stream.Failure ]
  | ("", "=") -> fun [ Equal -> "" | _ -> raise Stream.Failure ]
  | ("", s) ->
    fun
    [ Keyw k -> if k = s then "" else raise Stream.Failure
    | _ -> raise Stream.Failure ]
  | _ -> raise (Token.Error "bad token in match function") ]
;
```

### 9.10.2 The functor parameter

The type of the functor parameter is defined as:

```
module type GLexerType =
  sig
    type te = 'x;
    value lexer : Token.gllexer te;
  end;
```

The token type must be specified (type "te") and the lexer also, with the interface for lexers, of the gllexer type defined above, the record fields being described in the section "interface with the lexer", but with a general token type.

### 9.10.3 The resulting grammar module

Once a module of type "GLexerType" has been built (previous section, it is possible to create a grammar module by applying the functor "Grammar.GMake". For example:

```
module MyGram = Grammar.GMake MyLexer;
```

Notice that the function "Entry.parse" of this resulting module does not take a character stream as parameter, but a value of type "parsable". This function is equivalent to the function "parse\_parsable" of the non functorial interface. In short, the parsing of some character stream "cs" by some entry "e" of the example grammar above, must be done by:

```
MyGram.Entry.parse e (MyGram.parsable cs)
```

instead of:

```
MyGram.Entry.parse e cs
```

#### 9.10.4 GEXTEND and GDELETE\_RULE

The "GEXTEND" and "GDELETE\_RULE" statements are also added in the expressions of the ocaml language when the syntax extension kit "pa\_extend.cmo" is loaded. They have to be used for grammars defined with the functorial interface. Their syntax are:

```
expression ::= gextend
              | gdelete-rule
gdelete-rule ::= "GDELETE_RULE" gdelete-rule-body "END"
gextend      ::= "GEXTEND" gextend-body "END"
gextend-body ::= grammar-module-name extend-body
gdelete-rule-body ::= grammar-module-name delete-rule-body
grammar-module-name ::= qualid
```

See the syntax of the EXTEND statement for the meaning of the syntax entries not defined above.



# Chapter 10

## Pretty module

Pretty is a library module to pretty print data. It contains:

- The function "horiz\_vertic" to specify how data has to be printed.
- The function "sprintf" to format strings.
- The variable "line\_length" which is a reference specifying the maximum lines lengths.

### 10.1 Module description

#### 10.1.1 horiz\_vertic

The function "horiz\_vertic" takes two functions as parameters. When called, it calls its first function. If that function fails with some internal error that the function "sprintf" below may raise, the second function is called.

The type of "horiz\_vertic" is:

```
(unit -> 'a) -> (unit -> 'a) -> 'a
```

#### the horizontal function

The first function is said to be the "horizontal" function. It tries to pretty print the data on a single line. In the context of this function, if the strings built by the function "sprintf" (see below) contain newlines or have lengths greater than "line\_length", the function fails (with a internal exception local to the module).

#### the vertical function

In case of failure of the "horizontal function", the second function of "horiz\_vertic", the "vertical" function, is called. In the context of that function, the "sprintf" function behaves like the normal "sprintf" function of the OCaml library module "Printf".

#### 10.1.2 sprintf

The function "sprintf" works like its equivalent in the module "Printf" of the OCaml library, and takes the same parameters. Its difference is that if it is called in the context of the first function (the "horizontal" function) of the function "horiz\_vertic" (above), all strings built by "sprintf" are checked for newlines or length greater than the maximum line length. If it is the case, the "sprintf" function fails, and the horizontal function fails also.

If "sprintf" is not in the context of the horizontal function, it behaves like the usual "sprintf" function.

### 10.1.3 line\_length

The variable "line\_length" is a reference holding the maximum line length of lines printed horizontally. Its default is 77. This can be changed by the user before using "horiz\_vertic".

## 10.2 Example

Suppose you want to pretty print the XML code "<li>something</li>". If the "something" is short, you want to see:

```
<li>something</li>
```

If the "something" has several lines, you want to see that:

```
<li>
  something
</li>
```

A possible implementation is:

```
open Pretty;
horiz_vertic
  (fun () -> sprintf "<li>something</li>")
  (fun () -> sprintf "<li>\n  something\n</li>");
```

Notice that the "sprintf" above is the one of the library Pretty.

Notice also that, in a program displaying XML code, this "something" may contain other XML tags, and is therefore generally the result of other pretty printing functions, and the program should rather look like:

```
horiz_vertic
  (fun () -> sprintf "<li>%s</li>" (print something))
  (fun () -> sprintf "<li>\n  %s\n</li>" (print something))
```

Parts of this "something" can be printed horizontally and other vertically using other calls to "horiz\_vertic" in the user function "print" above. But it is important to remark that if they are called in the context of the first function parameter of "horiz\_vertic" above, only horizontal functions are accepted: the first failing "horizontal" function triggers the failure of the horizontal pretty printing.

## 10.3 Programming with Pretty

### 10.3.1 Hints

Just start with a call to "horiz\_vertic".

As its first function, use "sprintf" just to concat the strings without putting any newlines or indentations, just using e.g. spaces to separate pieces of data.

As its second function, wonder how you want your data to be cut. At the cutting point or points, add newlines. Notice that you probably need to give the current indentation string as parameter of the called functions because they need to be taken into account in the called "horizontal" functions.

In the example below, don't put the indentation in the sprintf function but give it as parameter of your "print" function:

```
horiz_vertic
  (fun () -> sprintf "<li>%s</li>" (print "" something))
  (fun () -> sprintf "<li>\n%s\n</li>" (print " " something))
```

Now, the "print" function could look like, supposing you print other things with "other" of the current indentation and "things" with a new shifted one:

```
value print ind something =
  horiz_vertic
    (fun () -> sprintf "%sother things..." ind)
    (fun () -> sprintf "%sother\n%s  things..." ind ind);
```

Supposing than "other" and "things" are the result of two other functions "print\_other" and "print\_things", your program could look like:

```
value print ind (x, y) =
  horiz_vertic
    (fun () -> sprintf "%s%s %s" ind (print_other 0 x) (print_things 0 y))
    (fun () -> sprintf "%s\n%s" (print_other ind x) (print_things (ind ^ " ") y));
```

### 10.3.2 How to cancel a horizontal print

If you want to prevent a pretty printing function to be called in an horizontal context, constraining the pretty print to be on several lines in the calling function, just do:

```
horiz_vertic
  (fun () -> sprintf "\n")
  (fun () -> ... (* your normal pretty print *))
```

In this case, the horizontal print always fails, due to the newline character in the sprintf format.

## 10.4 Remarks

### 10.4.1 Kernel

The module "Pretty" is supposed to be a basic, a "kernel" module to pretty print data. It supposes that the user takes care himself of the indentation. Programs using "Pretty" are not as short as the ones using "Format" of the OCaml library, but are more flexible. Later, it is planed to find a way to extend "Pretty" with functions allowing to use a short syntax similar to the "@" convention of the function "printf" of "Format", and taking care of the indentation for the user, resulting on shorter programs.

### 10.4.2 Strings vs Channels

In "Pretty", the pretty print is done only on strings, not on files. To pretty print on files, just built the strings and print them afterwards with the usual output functions. Notice that OCaml allocates and frees strings very fast, and if pretty printed values are not huge, which is generally the case, it is not a real problem, memory sizes these days being much more than enough for this job.

### 10.4.3 Strings or other types

The "horiz\_vertic" function can return values of other types than "string". For example, if you are interested only in the result of horizontal context and not on the vertical one, it is perfectly correct to write:

```
horiz_vertic
(fun () -> Some (sprintf "I hold on a single line"))
(fun () -> None)
```

### 10.4.4 Why raising exceptions ?

One could ask why this pretty print system has to raise internal exceptions. Why not simply write the pretty printing program like this:

1. first build the data horizontally (without newlines)
2. if the string length is lower than the maximum line length, return it
3. if not, build the string by adding newlines in the specific places

This method works but is generally very slow (exponential in time). Because while printing horizontally, many unuseful strings are built. If, for example, the final printed data holds on 50 lines, tenth of lines may be build and build again unusefully before the overflowing is tested.

# Chapter 11

## Printing programs

Camlp5 provides extensions kits to pretty print programs in revised syntax and normal syntax. Some other extensions kits also allow to rebuild the parsers, or the EXTEND statements in their initial syntax. The pretty print system is itself extensible, by adding new rules. We present here how it works in the Camlp5 sources.

The pretty print system of Camlp5 uses the library modules `Pretty`, an original system to format output) and `Extfun`, another original system of extensible functions.

This documentation is destinated to programmers who want to understand how the pretty printing of ocaml programs work in camlp5, want to adapt, modify or debug it, or want to add their own pretty printing extensions.

### 11.1 Introduction

The files doing the pretty prints are located in Camlp5 sources in the directory "etc". Look at them if you are interested on creating new ones. The main ones are:

- "etc/pr\_r.ml": pretty print in revised syntax.
- "etc/pr\_o.ml": pretty print in normal syntax.
- "etc/pr\_rp.ml": rebuilding parsers in their original revised syntax.
- "etc/pr\_op.ml": rebuilding parsers in their original normal syntax.
- "etc/pr\_extend.ml": rebuilding EXTEND in its original syntax.

We present here how this system work inside these files.

### 11.2 Principles

#### 11.2.1 Using module `Pretty`

All functions in ocaml pretty printing take a parameter named "the printing context" (variable `pc`). It is a record holding :

- The current indendation : `pc.ind`
- What has to be printed before, in the same line : `pc.bef`

- What has to be printed after, in the same line : `pc.aft`
- The dangling token, useful in normal syntax to know whether parentheses are necessary : `pc.dang`

A typical pretty printing function calls the function `horiz_vertic` of the library module `Pretty`. This function takes two functions as paramter:

- The way to print the data in one only line (*horizontal* printing)
- The way to print the data in two or more lines (*vertical* printing)

Both functions catenate the strings by using the function `sprintf` of the library module `Pretty` which controls whether the printed data holds in the line or not. They generally call, recursively, other pretty printing functions with the same behaviour.

Let us see an example (fictive) of printing an ocaml application. Let us suppose we have an application expression "`e1 e2`" to pretty print where `e1` and `e2` are sub-expressions. If both expressions and their application holds on one only line, we want to see:

```
e1 e2
```

On the other hand, if they do not hold on one only line, we want to see `e2` in another line with, say, an indendation of 2 spaces:

```
e1
  e2
```

Here is a possible implementation. The function has been named `expr_app` and can call the function `expr` to print the sub-expressions `e1` and `e2`:

```
value expr_app pc e1 e2 =
  horiz_vertic
    (fun () ->
      let s1 = expr {(pc) with aft = ""} e1 in
      let s2 = expr {(pc) with bef = ""} e2 in
      sprintf "%s %s" s1 s2)
    (fun () ->
      let s1 = expr {(pc) with aft = ""} e1 in
      let s2 =
        expr
          {(pc) with
            ind = pc.ind + 2;
            bef = tab (pc.ind + 2)}
          e2
      in
      sprintf "%s\n%s" s1 s2)
;
```

The first function is the *horizontal* printing. It ends with a `sprintf` separating the printing of `e1` and `e2` by a space. The possible "before part" (`pc.bef`) and "after part" (`pc.aft`) are transmitted in the calls of the sub-functions.

The second function is the *vertical* printing. It ends with a `sprintf` separating the printing of `e1` and `e2` by a newline. The second line starts with an indendation, using the "before part" (`pc.bef`) of the second call to `expr`.

The pretty printing library function `Pretty.horiz_vertic` calls the first (*horizontal*) function, and if it fails (either because `s1` or `s2` are too long or hold newlines, or because the final string produced by `sprintf` is too long), calls the second (*vertical*) function.

Notice that the parameter `pc` contains a field `pc.bef` (what has to be printed before in the same line), which in both cases is transmitted to the printing of `e1` (since the syntax `{(pc) with aft = ""}` is a record with `pc.bef` kept). Same for the field `pc.aft` transmitted to the printing of `e2`.

## 11.2.2 Using module Extfun and its syntax

This system is combined to the extensible functions to allow the extensibility of the pretty printing. Pretty printers of `camlp5` can then be used as "kits" to be added or not, according to the things to be pretty printed in some or other ways. In particular, the pretty printing kit `"pr_r.cmo"` alone does not rebuild parsers in their original syntax. When adding `"pr_rp.cmo"`, the parsers are rebuilt: the code of `"pr_rp.ml"` is just an extension of some parts of the pretty printing extensible functions of `"pr_r.ml"`.

The code above actually looks like:

```
value expr_app =
  extfun Extfun.empty with
  [ <:expr< $e1$ $e2$ >> ->
    fun curr next pc ->
      horiz_vertic
        (fun () ->
          let s1 = curr {(pc) with aft = ""} e1 in
          let s2 = next {(pc) with bef = ""} e2 in
          sprintf "%s %s" s1 s2)
        (fun () ->
          let s1 = curr {(pc) with aft = ""} e1 in
          let s2 =
            next
              {(pc) with
                ind = pc.ind + 2;
                bef = tab (pc.ind + 2)}
              e2
          in
          sprintf "%s\n%s" s1 s2)
    | e ->
      fun curr next pc -> next pc e ]
;
```

The extensible functions have a syntax tree (here `<:expr< $e1$ $e2$ >>`) as parameter. To be extensible, the syntax tree must be the first parameter (it is not possible to apply extensions inside a closure). The other parameters, in particular the printing context `pc` are given in the semantic action.

The parameter `curr` and `next` are provided by the pretty printing system for ocaml programs. They correspond to the pretty printing of, respectively, the current level and the next level. Since the application in ocaml is left associative, the first sub-expression is printed at the same (current) level and the second one is printed at the next level. We also see a call to `next` in the last (2nd) case of the function to treat the other cases in the next level.

### 11.2.3 Dangling else, bar, semicolon

In normal syntax, there are cases where it is necessary to enclose expressions between parentheses (or between begin and end, which is equivalent in that syntax). Three tokens may cause problems: the "else", the vertical bar "|" and the semicolon ";". Here are examples where the presence of these tokens constraints the previous expression to be parenthesized. In these three examples, removing the begin..end enclosers would change the meaning of the expression because the dangling token would be included in that expression:

Dangling else:

```
if a then begin if b then c end else d
```

Dangling bar:

```
function
  A ->
    begin match a with
      B -> c
      | D -> e
    end
  | F -> g
```

Dangling semicolon:

```
if a then b
else begin
  let c = d in
  e
end;
f
```

The information is transmitted by the value `pc.dang`. In the first example above, while displaying the "then" part of the outer "if", the sub-expression is called with the value `pc.dang` set to "else" to inform the last sub-sub-expression that it is going to be followed by that token. When a "if" expression has to be displayed without "else" part, and that its "`pc.dang`" is "else", it has to be enclosed with spaces.

This problem does not exist in revised syntax. While pretty printing in revised syntax, the parameter `pc.dang` is not necessary and remains the empty string.

### 11.2.4 By level

For each level of pretty printing, there is such a function. The example showed the pretty printing of expression at the level "apply". There are other functions for levels "top", "add", "mul", "simple", and so on. The global pretty printing variable for expressions is a record, named "`pr_expr`" (in the module `Pcaml.Printers`), where the levels are defined by a list, something like this:

```
pr_expr.pr_levels :=
  [{pr_label = "top"; pr_rules = expr_top};
   {pr_label = "add"; pr_rules = expr_add};
   {pr_label = "mul"; pr_rules = expr_mul};
   {pr_label = "apply"; pr_rules = expr_app};
   {pr_label = "simple"; pr_rules = expr_simple}]
;
```

where we find, in particular, our function `expr_app` defined above.

The call to a specific level is done by the function `pr_expr.pr_fun` with the level name. It returns the function taking the "printing context" (`pc`) and the expression as parameters, and returning the pretty printed string. For example, the call to the top level of expressions has been defined as:

```
value expr pc e = pr_expr.pr_fun "top" pc e;
```

Same thing for the other pretty printed functions for patterns, structures, signatures, and so on.

To extend some level, in another file, the function `find_pr_level` can be used to get the level to be extended, e.g.

```
let expr_app = find_pr_level "app" pr_expr.pr_levels in
expr_app.pr_rules :=
  extfun expr_app.pr_rules with
  [ <:expr< .... >> -> ...
  | <:expr< .... >> -> ...
  | ... ];
```



# Chapter 12

## Scheme and Lisp syntaxes

It is possible to write ocaml programs with Scheme or Lisp syntax. They are close the one to the other, both using parentheses to identify and separate statements.

### 12.1 Common

The syntax extension kits are named `"pa_scheme.cmo"` and `"pa_lisp.cmo"`. The sources (same names ending with `".ml"` in the `camlp5` sources), are written in their own syntax. They are bootstrapped thanks to versions written in revised syntax and to the `camlp5` pretty printing system.

In the ocaml toplevel, it is possible to use them by loading `"camlp5r.cma"` first, then `"pa_lisp.cmo"` or `"pa_scheme.cmo"` after:

```
ocaml -I +camlp5 camlp5r.cma pa_scheme.cmo
      Objective Caml version ...
```

```
      Camlp5 Parsing version ...
```

```
# (let ((x 3)) (* 3 x))
- : int = 9
# (values 3 4 5)
- : (int * int * int) = (3, 4, 5)
```

```
ocaml -I +camlp5 camlp5r.cma pa_lisp.cmo
      Objective Caml version ...
```

```
      Camlp5 Parsing version ...
```

```
# (let ((x 3)) (* 3 x))
- : int = 9
# (, 3 4 5)
- : (int * int * int) = (3, 4, 5)
```

The grammar of Scheme and Lisp are relatively simple, just reading s-expressions. The syntax tree nodes are created in the semantic actions. Because of this, these grammars are hardly extensible.

However, the syntax extension EXTEND ("pa\_extend.cmo" in extensible grammars) works for them: only the semantic actions have to be written with the Scheme or Lisp syntax. The stream parsers are also implemented.

Warning: these syntaxes are incomplete, but can be completed, if camlp5 users are interested.

## 12.2 Scheme syntax

Some examples are given to show the principles:

OCaml	Scheme
<pre> let x = 42;; let f x = 0;; let rec f x y = 0;; let x = 42 and y = 27 in x + y;; let x = 42 in let y = 27 in x + y;; module M = struct ... end;; type 'a t = A of 'a * int   B fun x y -&gt; x x; y; z f x y [1; 2; 3] x :: y :: z :: t a, b, c match x with 'A'..'Z' -&gt; "x" {x = y; z = t} </pre>	<pre> (define x 42) (define (f x) 0) (definerec (f x y) 0) (let ((x 42) (y 27)) (+ x y)) (let* ((x 42) (y 27)) (+ x y)) (module M (struct ...)) (type (t 'a) (sum (A 'a int) (B))) (lambda (x y) x) (begin x y z) (f x y) [1 2 3] [x y z :: t] (values a b c) (match x ((range 'A' 'Z') "x")) {(x y) (z t)} </pre>

## 12.3 Lisp syntax

The same examples:

OCaml	Lisp
<pre> let x = 42;; let f x = 0;; let rec f x y = 0;; let x = 42 and y = 27 in x + y;; let x = 42 in let y = 27 in x + y;; module M = struct ... end;; type 'a t = A of 'a * int   B fun x y -&gt; x x; y; z f x y [1; 2; 3] x :: y :: z :: t a, b, c match x with 'A'..'Z' -&gt; "x" {x = y; z = t} </pre>	<pre> (value x 42) (value f (lambda x 0)) (value rec f (lambda (x y) 0)) (let ((x 42) (y 27)) (+ x y)) (let* ((x 42) (y 27)) (+ x y)) (module M (struct ...)) (type (t 'a) (sum (A 'a int) (B))) (lambda (x y) x) (progn x y z) (f x y) (list 1 2 3) (list x y z :: t) (, a b c) (match x ((range 'A' 'Z') "x")) ({} (x y) (z t)) </pre>

# Chapter 13

## Syntax tree

In Camlp5, one often uses syntax trees. For example, in grammars of the language (semantic actions), in pretty printing (as patterns), in optimizing syntax code (typically streams parsers). Syntax trees are mainly defined by sum types, one for each kind of tree: `"expr"` for expressions, `"patt"` for patterns, `"ctyp"` for types, `"str_item"` for structure items, and so on. Each node corresponds to a possible value of this type.

### 13.1 Introduction

This syntax tree is defined in the module `"MLast"` provided by Camlp5.

For example, the syntax tree of the statement `"if"` can be written:

```
MLast.ExIfe loc e1 e2 e3
```

where `"loc"` is the location in the source, and `"e1"`, `"e2"` and `"e3"` are respectively the expression after the `"if"`, the one after the `"then"` and the one after the `"else"`.

In all programs, it is possible to manipulate syntax trees like that.

However, the quotations systems of Camlp5 provides a quotation kit named `"q_MLast.cmo"`. When loaded, it is possible to represent the syntax trees in concrete syntax. The example above can be written:

```
<:expr< if $e1$ then $e2$ else $e3$ >>
```

The interest of this representation is when one must manipulate complicated syntax trees. An example taken from the Camlp5 sources is this quotation (found in a pattern):

```
<:expr<
  match try Some $f$ with [ Stream.Failure -> None ] with
  [ Some $p$ -> $e$
    | _ -> raise (Stream.Error $e2$) ]
>>
```

In abstract syntax, it should have been written:

```
MLast.ExMat _
  (MLast.ExTry _ (MLast.ExApp _ (MLast.ExUId _ "Some") f)
    [(MLast.PaAcc _ (MLast.PaUId _ "Stream") (MLast.PaUId _ "Failure"),
      None, MLast.ExUId _ "None"]])
```

```

    [(MLast.PaApp _ (MLast.PaUId _ "Some") p, None, e);
     (MLast.PaAny _, None,
      MLast.ExApp _ (MLast.ExLid _ "raise")
        (MLast.ExApp _
          (MLast.ExAcc _ (MLast.ExUId _ "Stream") (MLast.ExUId _ "Error"))
          e2))]

```

Which is less readable.

Instead of thinking of "a syntax tree", the programmer has to think of "a piece of program".

## 13.2 Location

In all syntax tree nodes, the first parameter is the source location of the node.

### 13.2.1 In expressions

When a quotation is in the context of an expression, the location parameter is "loc" in the node and in all its possible sub-nodes. Example: if we consider the quotation:

```
<:sig_item< value foo : int -> bool >>
```

This quotation, in a context of an expression, is equivalent to:

```

MLast.SgVal loc "foo"
  (MLast.TyArr loc (MLast.TyLid loc "int") (MLast.TyLid loc "bool"));

```

This name, "loc", is predefined, but it is possible to change it, using the argument "-loc" of the `camp5` shell commands.

Consequently, if there is no variable "loc" defined in the context of the quotation, or if it is not of the good type, a semantic error occur in the `ocaml` compiler.

Note that in the extensible grammars, the variable "loc" is bound, in all semantic actions, to the location of the rule.

If the created node has no location, the simplest way is to define the variable "loc" to have the value "Stdpp.dummy\_loc".

### 13.2.2 In patterns

When a quotation is in the context of a pattern, the location parameter of all nodes and possible sub-nodes is set to the wildcard ("\_"). The same example above:

```
<:sig_item< value foo : int -> bool >>
```

is equivalent, in a pattern, to:

```

MLast.SgVal _ "foo"
  (MLast.TyArr _ (MLast.TyLid _ "int") (MLast.TyLid _ "int"))

```

## 13.3 Antiquotations

The expressions or patterns between dollar (\$) characters are called *antiquotations*. In opposition to quotations which has its own syntax rules, the antiquotation is an area in the syntax of the enclosing context (expression or pattern). See the chapter about quotations.

If the quotation (any quotation) is in the context of an expression, the antiquotation is an expression. It could be a call to a function. Examples:

```
value f e el = <:expr< [$e$ :: $loop False el$] >>;
value patt_list p pl = <:patt< ( $list:[p::pl]$) >>;
```

If the quotation (any quotation) is in the context of a pattern, the antiquotation is a pattern. Any pattern is possible, including the wildcard character ("\_"). Examples:

```
fun [ <:expr< $lid:op$ $_$ $_$ >> -> op ]
match p with [ <:patt< $_$ | $_$ >> -> Some p ]
```

## 13.4 Nodes and Quotations

This section describes all nodes defined in the module "MLast" of camlp5 and how to write them with quotations. Notice that, inside quotations, one is not restricted to these elementary cases, but any complex value can be used, resulting on possibly complex combined nodes.

Variables names give information of their types:

- e, e1, e2, e3: expr
- p, p1, p2, p3: patt
- t, t1, t2, e3: ctyp
- s: string
- b: bool
- me, me1, me2: module\_expr
- mt, mt1, mt2: module\_type
- le: list expr
- lp: list patt
- lt: list ctyp
- ls: list string
- lse: list (string \* expr)
- lpe: list (patt \* expr)
- lpp: list (patt \* patt)
- lpoe: list (patt \* option expr \* expr)
- op: option patt
- lcstri: list class\_str\_item
- lcsigi: list class\_sig\_item

### 13.4.1 expr

Expressions of the language.

Node	<:expr< ... >>	Comment
ExAcc loc e1 e2	\$e1\$ . \$e2\$	dot
ExAnt loc e	\$ante:e\$	antiquotation (1)
ExApp loc e1 e2	\$e1\$ \$e2\$	application
ExAre loc e1 e2	\$e1\$ .( \$e2\$ )	array access
ExArr loc le	[  \$list:le\$  ]	array
ExAsr loc e	assert \$e\$	assert
ExAss loc e1 e2	\$e1\$ := \$e2\$	assignment
ExBae loc e le	\$e\$ .{ \$le\$ }	big array access
ExChr loc s	\$chr:s\$	character constant
ExCoe loc e (Some t1) t2	(\$e1\$ : \$t1\$ :> \$t2\$)	coercion
ExCoe loc e None t2	(\$e1\$ :> \$t2\$)	coercion
ExFlo loc s	\$flo:s\$	float constant
ExFor loc s e1 e2 b le	for \$\$ = \$e1\$ \$to:b\$ \$e2\$ do { \$list:le\$ }	for
ExFun loc lpoee	fun [ \$list:lpoee\$ ]	function (2)
ExIfe loc e1 e2 e3	if \$e1\$ then \$e2\$ else \$e3\$	if
ExInt loc s ""	\$int:s\$	integer constant
ExInt loc s "1"	\$int32:s\$	integer 32 bits
ExInt loc s "L"	\$int64:s\$	integer 64 bits
ExInt loc s "n"	\$nativeint:s\$	native integer
ExLab loc s None	\$\$	label
ExLab loc s (Some e)	\$\$ : \$e\$	label
ExLaz loc e	lazy \$e\$	lazy
ExLet loc b lpe e	let \$opt:b\$ \$list:lpe\$ in \$e\$	let binding
ExLid loc s	\$lid:s\$	lowercase identifier
ExLmd loc s me e	let module \$\$ = \$me\$ in \$e\$	let module
ExMat loc e lpoe	match \$e\$ with [ \$list:lpoe\$ ]	match (2)
ExNew loc ls	new \$list:ls\$	new
ExObj loc op lcstri	object (\$opt:op\$) \$list:lcstri\$ end	object expression
ExOlb loc s None	? \$\$	option label
ExOlb loc s (Some e)	? \$\$ : \$e\$	option label
ExOvr loc lse	{< \$lse\$ >}	override
ExRec loc lpe None	{ \$list:lpe\$ }	record
ExRec loc lpe (Some e)	{ (\$e\$) with \$list:lpe\$ }	record
ExSeq loc le	do { \$list:le\$ }	sequence
ExSnd loc e s	\$e\$ # \$\$	method call
ExSte loc e1 e2	\$e1\$ .[ \$e2\$ ]	string element
ExStr loc s	\$str:s\$	string
ExTry loc e lpoe	try \$e\$ with [ \$list:lpoe\$ ]	try (2)
ExTup loc le	(\$list:le\$)	t-uple
ExTyc loc e t	(\$e\$ : \$t\$)	type constraint
ExUid loc s	\$uid:s\$	uppercase identifier
ExVrn loc s	' \$\$	variant
ExWhi loc e le	while \$e\$ do { \$list:le\$ }	while

(1) Node used to specify an antiquotation area. See the chapter about quotations.

(2) The variable "lpoe" found in "function", "match" and "try" statements correspond to a list of "(patt \* option expr \* expr)" where the "option expr" is the "when" optionally following the pattern:

p -> e

is represented by:

(p, None, e)

and

p when e1 -> e

is represented by:

(p, Some e1, e)

### 13.4.2 patt

Patterns of the language.

Node	<:patt< ... >>	Comment
PaAcc loc p1 p2	\$p1\$ . \$p2\$	dot
PaAli loc p1 p2	(\$p1\$ as \$p2\$)	alias
PaAnt loc p	\$anti:e\$	antiquotation (1)
PaAny loc	-	wildcard
PaApp loc p1 p2	\$p1\$ \$p2\$	application
PaArr loc lp	[ \$list:lp\$  ]	array
PaChr loc s	\$chr:s\$	character
PaInt loc s1 s2	\$int:s\$	integer
PaFlo loc s	\$flo:s\$	float
PaLab loc s None	\$s\$	label
PaLab loc s (Some p)	\$s\$ : \$p\$	label
PaLid loc s	\$lid:s\$	lowercase identifier
PaOlb loc s None	? \$s\$	option label
PaOlb loc s (Some (p, None))	? \$s\$ : (\$p\$)	option label
PaOlb loc s (Some (p, Some e))	? \$s\$ : (\$p\$ = \$e\$)	option label
PaOrp loc p1 p2	\$p1\$   \$p2\$	or
PaRng loc p1 p2	\$p1\$ .. \$p2\$	range
PaRec loc lpp None	{ \$list:lpp\$ }	record
PaStr loc s	\$str:s\$	string
PaTup loc lp	(\$list:lp\$)	t-uple
PaTyc loc p t	(\$p\$ : \$t\$)	type constraint
PaTyp loc ls	# \$list:ls\$	type pattern
PaUid loc s	\$uid:s\$	uppercase identifier
PaVrn loc s	' \$s\$	variant

(1) Node used to specify an antiquotation area. See the chapter about quotations.

### 13.4.3 ctyp

Type expressions of the language.

Node	<:ctyp< ... >>	Comment
TyAcc loc t1 t2	\$t1\$ . \$t2\$	dot
TyAli loc t1 t2	\$t1\$ as \$t2\$	alias
TyAny loc	-	wildcard
TyApp loc t1 t2	\$t1\$ \$t2\$	application
TyArr loc t1 t2	\$t1\$ -> \$t2\$	arrow
TyCls loc ls	# \$list:ls\$	class
TyLab loc s t	\$s\$ : \$t\$	label
TyLid loc s	\$lid:s\$	lowercase identifier
TyMan loc t1 t2	\$t1\$ == \$t2\$	manifest
TyObj loc lst False	< \$list:lst\$ >	object
TyObj loc lst True	< \$list:lst\$ .. >	object
TyObj loc lst b	< \$list:lst\$ \$opt:b\$ >	object (general)
TyOlb loc s t	? \$s\$ : \$t\$	option label
TyPol loc ls t	! \$list:ls\$ . \$t\$	polymorph
TyQuo loc s	' \$s\$	variable
TyRec loc llsbt	{ \$list:llsbt\$ }	record
TySum loc llslt	[ \$list:llslt\$ ]	sum
TyTup loc lt	( \$list:lt\$ )	t-uple
TyUid loc s	\$uid:s\$	uppercase identifier
TyVrn loc lpv None	[ = \$list:lpv\$ ]	variant
TyVrn loc lpv (Some None)	[ > \$list:lpv\$ ]	variant
TyVrn loc lpv (Some (Some []))	[ < \$list:lpv\$ ]	variant
TyVrn loc lpv (Some (Some ls))	[ < \$list:lpv\$ > \$list:ls\$ ]	variant

### 13.4.4 modules...

#### str\_item

Structure items, i.e. phrases in a ".ml" file or "struct"s elements.

Node	<:str_item< ... >>	Comment
StCls loc lcd	class \$list:lcd\$	class declaration
StClT loc lcdt	class type \$list:lctd\$	class type declaration
StDcl loc lstri	declare \$list:lstri\$ end	declare
StDir loc s None	# \$s\$	directive
StDir loc s (Some e)	# \$s\$ \$e\$	directive
StDir loc s oe	# \$s\$ \$opt:oe\$	directive (general)
StExc loc s lt []	exception \$s\$ of \$list:lt\$	exception
StExc loc s lt ls	exception \$s\$ of \$list:lt\$ = \$list:ls\$	exception
StExp loc e	\$exp:e\$	expression
StExt loc s t ls	external \$s\$ : \$t\$ = \$list:ls\$	external
StInc loc me	include \$me\$	include
StMod loc b lsme	module \$opt:b\$ \$list:lsme\$	module
StMty loc s mt	module type \$s\$ = \$mt\$	module type

StOpn loc ls	open \$list:ls\$	open
StTyp loc ltd	type \$list:ltd\$	type declaration
StUse loc s lstrib	...internal use...	
StVal loc b lpe	value \$opt:b\$ \$list:lpe\$	value

## sig\_item

Signature items, i.e. phrases in a ".mli" file or "sig"s elements.

Node	<:sig_item< ... >>	Comment
SgCls loc lcd	class \$list:lcd\$	class
SgClt loc lct	class type \$list:lct\$	class type
SgDcl loc lsigi	declare \$list:lsigi\$ end	declare
SgDir loc s None	# \$\$	directive
SgDir loc s (Some e)	# \$\$\$ \$e\$	directive
SgDir loc s oe	# \$\$\$ \$opt:oe\$	directive (general)
SgExc loc s []	exception \$\$	exception
SgExc loc s lt	exception \$\$ of \$list:lt\$	exception
SgExt loc s t ls	external \$\$ : \$t\$ = \$list:ls\$	external
SgInc loc me	include \$me\$	include
SgMod loc b lsmt	module \$opt:b\$ \$list:lsmt\$	module
SgMty loc s mt	module type \$\$ = \$mt\$	module type
SgOpn loc ls	open \$list:ls\$	open
SgTyp loc ltd	type \$list:ltd\$	type declaration
SgUse loc s lstrib	...internal use...	
SgVal loc s t	value \$\$ : \$t\$	value

## module\_expr

Node	<:module_expr< ... >>	Comment
MeAcc loc me1 me2	\$me1\$ . \$me2\$	dot
MeApp loc me1 me2	\$me1\$ \$me2\$	application
MeFun loc s mt me	functor ( \$\$ : \$mt\$ ) -> \$me\$	functor
MeStr loc lstri	struct \$list:lstri\$ end	struct
MeTyc loc me mt	( \$me\$ : \$mt\$ )	module type constraint
MeUid loc s	\$uid:s\$	uppercase identifier

## module\_type

Node	<:module_type< ... >>	Comment
MtAcc loc mt1 mt2	\$mt1\$ . \$mt2\$	dot
MtApp loc mt1 mt2	\$mt1\$ \$mt2\$	application
MtFun loc s mt1 mt2	functor ( \$\$ : \$mt1\$ ) -> \$mt2\$	functor

MtLid loc s	\$lid:s\$	lowercase identifier
MtQuo loc s	' \$\$	abstract
MtSig loc lsigi	sig \$list:lsigi\$ end	signature
MtUId loc s	\$uid:s\$	uppercase identifier
MtWit loc mt lwc	\$mt\$ with \$list:lwc\$	with construction

### 13.4.5 classes...

#### class\_expr

Node	<:class_expr< ... >>	Comment
CeApp loc ce e	\$ce\$ \$e\$	application
CeCon loc ls lt	\$list:ls\$ [ \$list:lt\$ ]	constructor
CeFun loc p ce	fun \$p\$ -> \$ce\$	function
CeLet loc b lpe ce	let \$opt:b\$ \$list:lpe\$ in \$ce\$	let binding
CeStr loc po lcstri	object (\$opt:op\$) \$list:lcstri\$ end	object
CeTyc loc ce ct	(\$ce\$ : \$ct\$)	class type constraint

#### class\_type

Node	<:class_type< ... >>	Comment
CtCon loc ls lt	\$list:ls\$ [ \$list:lt\$ ]	constructor
CtFun loc t ct	[ \$t\$ ] -> \$ct\$	arrow
CtSig loc pt None lcsigi_item	object \$list:lcsigi\$ end	object
CtSig loc pt (Some t) lcsigi_item	object (\$t\$) \$list:lcsigi\$ end	object
CtSig loc pt ot lcsigi_item	object \$opt:ot\$ \$list:lcsigi\$ end	object (general)

#### class\_str\_item

Node	<:class_str_item< ... >>	Comment
CrCtr loc t1 t2	type \$t1\$ = \$t2\$	type constraint
CrDcl loc lcstri	declare \$list:lcstri\$ end	declaration list
CrInh loc ce None	inherit \$ce\$	inheritance
CrInh loc ce (Some s)	inherit \$ce\$ as \$s\$	inheritance
CrInh loc ce os	inherit \$ce\$ \$opt:s\$	inheritance (general)
CrIni loc e	initializer \$e\$	initialization
CrMth loc s False e None	method \$\$ = \$e\$	method
CrMth loc s False e (Some t)	method \$\$ : \$t\$ = \$e\$	method
CrMth loc s True e None	method private \$\$ = \$e\$	method
CrMth loc s True e (Some t)	method private \$\$ : \$t\$ = \$e\$	method
CrMth loc s b e ot	method \$opt:b\$ \$\$ \$opt:ot\$ = \$e\$	method (general)
CrVal loc s False e	value \$\$ = \$e\$	value
CrVal loc s True e	value mutable \$\$ = \$e\$	value
CrVal loc s b e	value \$opt:b\$ \$\$ = \$e\$	value (general)

CrVir loc s False t	method virtual \$\$ : \$t\$	virtual method
CrVir loc s True t	method virtual private \$\$ : \$t\$	virtual method
CrVir loc s b t	method virtual \$opt:b\$ \$\$ : \$t\$	virtual method (general)

### class\_sig\_item

Node	<:class_sig_item< ... >>	Comment
CgCtr loc t1 t2	type \$t1\$ = \$t2\$	type constraint
CgDcl loc s lcsigi	declare \$list:lcsigi\$ end	declare
CgInh loc ct	inherit \$ct\$	inheritance
CgMth loc s False t	method \$\$ : \$t\$	method
CgMth loc s True t	method private \$\$ : \$t\$	method
CgMth loc s b t	method \$opt:b\$ \$\$ : \$t\$	method (general)
CgVal loc s False t	value \$\$ : \$t\$	value
CgVal loc s True t	value mutable \$\$ : \$t\$	value
CgVal loc s b t	value \$opt:b\$ \$\$ : \$t\$	value (general)
CgVir loc s False t	method virtual \$\$ : \$t\$	method virtual
CgVir loc s True t	method virtual private \$\$ : \$t\$	method virtual
CgVir loc s b t	method virtual \$opt:b\$ \$\$ : \$t\$	method virtual (general)

## 13.4.6 other

### with\_constr

”With” possibly following a module type.

Node	<:with_const< ... >>	Comment
WcTyp loc s ltv False t	type \$\$ \$list:ltv\$ = \$t\$	with type
WcTyp loc s ltv True t	type \$\$ \$list:ltv\$ = private \$t\$	with type
WcTyp loc s ltv b t	type \$\$ \$list:ltv\$ = \$opt:b\$ \$t\$	with type (general)
WcMod loc ls me	module \$list:ls\$ = \$me\$	with module

### poly\_variant

Polymorphic variants.

Node	<:poly_variant< ... >>	Comment
PvTag s False []	‘ \$i\$	constructor
PvTag s True lt	‘ \$i\$ of & \$list:lt\$	constructor
PvTag s b lt	‘ \$i\$ of \$opt:b\$ \$list:lt\$	constructor (general)
PvInh t	\$t\$	type



# Chapter 14

## Locations

The location is a concept often used in `camlp5`, bound to know where the errors occur in the source. The basic type is `"Stdpp.location"` which is an abstract type.

### 14.1 Definitions

A location is internally a couple of source *positions*: the beginning and the end of an element in the source (file or interactive). A located element can be a character (the end is just the beginning plus one), a token, or a longer sequence generally corresponding to a grammar rule.

A *position* is a count of characters since the beginning of the file, starting at zero. When a couple of positions define a location, the first position is the position of the first character of the element, and the last position is the first character *not* part of the element. The location length is the difference between those two numbers. Notice that the position corresponds exactly to the character count in the streams of characters.

In the extensible grammars, a variable with the specific name `"loc"` is predefined in all semantic actions: it is the location of the associated rule. Since the syntax tree quotations generate nodes with `"loc"` as location part, this allow to generate grammars without having to think about source locations.

It is possible to change the name `"loc"` into another name, through the parameter `"-loc"` of the `camlp5` commands.

Remark: the reason why the type `"location"` is abstract is that in future versions, it may contain other informations, such as the associated comments, the type (for expressions nodes), things like that, without having to change the already written programs.

### 14.2 Building locations

Tools are provided in the module `"Stdpp"` to manage locations.

First, `"Stdpp.dummy_loc"` is a dummy location used when the element does not correspond to any source, or if the programmer does not want to busy about locations.

The function `"Stdpp.make_lined_loc"` builds a location from three parameters:

- the line number, starting at 1

- the position of the first column of the line
- a couple of positions of the location: the first one belonging to the given line, the second one being able to belong to another line, further.

If the line number is not known, it is possible to use the function `"Stdpp.make_loc"` taking only the couple of positions of the location. In this case, error messages may indicate the first line and a big count of characters from this line (actually from the beginning of the file). With a good text editor, it is possible, anyway to find the good location, anyway.

If the location is built with `"Stdpp.make_loc"`, and if your program displays a source location itself, it is possible to use the function `"Stdpp.line_of_loc"` which takes the file name and the location as parameters and return, by reading that file, the line number, and the character positions of the location.

## 14.3 Raising with a location

The function `"Stdpp.raise_with_loc"` allows to raise an exception together with a location. It is the case of all exceptions raised in the extensible grammars. The raised exception is `"Stdpp.Exc_located"` with two parameters: the location and the exception itself.

Notice that `"raise_with_loc"` just reraises the exception if it is already the exception `"Exc_located"`, ignoring then the new given location.

A good usage to print exceptions possibly enclosed by `"Exc_located"` is to write the `"try..with"` statement like this:

```
try ... with exn ->
  let exn =
    match exn with
    [ Stdpp.Exc_located loc exn -> do { ... print the location ...; exn }
    | _ -> exn ]
  in
  match exn with
  ...print the exception which is *not* located...
```

## 14.4 Other functions

Some other functions are provided:

`Stdpp.first_pos`

returns the first position (an integer) of the location.

`Stdpp.last_pos`

returns the last position (an integer) of the location (position of the first character not belonging to the element).

`Stdpp.line_nb`

returns the line number of the location or -1 if the location does not contain a line number (i.e. built by `"Stdpp.make_loc"`).

#### `Stdpp.bol_pos`

returns the position of the beginning of the line of the location. It is zero if the location does not contain a line number (i.e. built by `"Stdpp.make_loc"`).

And still other ones used in camlp5 sources:

#### `Stdpp.encl_loc`

`"Stdpp.encl_loc loc1 loc2"` returns the location starting at the smallest begin of `"loc1"` and `"loc2"` and ending at their greatest end.. In simple words, it is the location enclosing `"loc1"` and `"loc2"` and all what is between them.

#### `Stdpp.shift_loc`

`"shift_loc sh loc"` returns the location `"loc"` shifted with `"sh"` characters. The line number is not recomputed.

#### `Stdpp.sub_loc`

`"Stdpp.sub_loc loc sh len"` is the location `"loc"` shifted with `"sh"` characters and with length `"len"`. The previous ending position of the location is lost.

#### `"Stdpp.after_loc"`

`"Stdpp.after_loc loc sh len"` is the location just after `"loc"` (i.e. starting at the end position of `"loc"`), shifted with `"sh"` characters, and of length `"len"`.



## Chapter 15

# Pragma directive

The directive `#pragma` allows to evaluate expressions at parse time, useful, for example, to test syntax extensions by the statement `EXTEND` without having to compile it in a separate file.

To use it, add the syntax extension `pa_pragma.cmo` in the `camlp5` command line. It add the ability to use this directive.

Example in a file, adding a syntax for the statement `'repeat'` and using it just after:

```
#pragma
EXTEND
  GLOBAL: expr;
  expr: LEVEL "top"
    [ [ "repeat"; e1 = sequence; "until"; e2 = SELF ->
      <:expr< do { $e1$; while not $e2$ do { $e1$ } } >> ] ]
    ;
  sequence:
    [ [ e1 = LIST1 expr_semi -> <:expr< do { $list:e1$ } >> ] ]
    ;
  expr_semi:
    [ [ e = expr; ";" -> e ] ]
    ;
END;

let i = ref 1 in
repeat print_int i.val; print_endline ""; incr i; until i.val = 10;
```

The compilation of this example (naming it `"foo.ml"`) can be done with the command:

```
ocamlc -pp "camlp5r q_MLast.cmo pa_extend.cmo pa_pragma.cmo" -I +camlp5 foo.ml
```

Notice that it is still experimental and probably incomplete, for the moment.



# Chapter 16

## Extensible functions

Extensible functions allow to define functions by pattern matching which are extensible by adding of new cases which are inserted automatically at the good place by comparing the patterns. The pattern cases are ordered according to syntax trees representing them, "when" statements being inserted before the cases without "when".

Notice that extensible functions are functional: when extending a function, a new function is returned.

The extensible functions are used in the pretty printing system of Camlp5.

### 16.1 Syntax

The syntax of the extensible functions, when loading "pa\_extfun.cmo" is the following:

```
expression ::= extensible-function
extensible-function ::= "extfun" expression "with" "[" match-cases "]"
match-cases ::= match-case "|" match-cases
match-case ::= pattern "->" expression
               | pattern "when" expression "->" expression
```

It is actually the same syntax than the one of "match" and "try" constructions.

### 16.2 Semantics

The statement "extend" defined by the syntax takes an extensible function and return another extensible function with the new match cases inserted at good places into the initial extensible function.

Extensible functions are of type "Extfun.t a b", which is an abstract type, where "a" and "b" are respectively the type of the patterns and the type of the expressions. It corresponds to a function of type "a -> b".

The function "Extfun.apply" takes an extensible function as parameter and return the function which can be applied like a normal function.

The value "Extfun.empty" is an empty extensible function, of type "Extfun.t 'a 'b". When applied with "Extfun.apply" and a parameter, it raises the exception "Extfun.Failure" whatever the parameter.

For debugging, it is possible to use the function `"Extfun.print"` which displays the match cases of the extensible functions. Actually, only the patterns are displayed in clear, the associated expressions are not.

The match cases are inserted according to the following rules:

- The match cases are inserted in the order they are defined in the syntax `"extend"`
- A match case pattern with `"when"` is inserted before a match case pattern without `"when"`.
- Two match cases patterns both with `"when"` or both without `"when"` are inserted according to the alphabetic order of some internal syntax tree of the patterns where bound variables names are not taken into account.
- If two match cases patterns without `"when"` have the same patterns internal syntax tree, the initial one is silently removed.
- If two match cases patterns with `"when"` have the same patterns internal syntax tree, the new one is inserted before the old one.