

Разработка программ с помощью Objective
Caml
Developing Applications with Objective Caml

Emmanuel Chailloux
Pascal Manoury
Bruno Pagano

26 августа 2012 г.

Оглавление

1	Как заполнить Objective CAML	5
I	Основы языка	7
2	Функциональное программирование	9
3	Императивное программирование	11
4	Функциональный и императивный стиль	13
5	Графический интерфейс	15
6	Приложения	17
II	Средства разработки	19
7	Компиляция и переносимость	21
8	Библиотеки	23
9	Автоматический сборщик мусора	25
10	Средства анализа программ	27
11	Средства лексического и синтаксического анализа	29
11.1	Введение	29
11.2	План главы	30
11.3	Лексика	30
11.3.1	Модуль Genlex	31
11.3.2	Использование потоков	32
11.4	Регулярные выражения	33

11.4.1 Инструмент Ocamllex	36
12 Взаимодействие с языком C	39
13 Приложения	41
 III Устройство программы	 43
14 Модульное программирование	45
15 Объектно-ориентированное программирование	47
16 Сравнение моделей устройств программ	49
17 Приложения	51
 IV Параллелизм и распределение	 53
18 Процессы и связь между процессами	55
18.1 Параллельное программирование	55
19 Распределённое программирование	57
20 Приложения	59
 V Разработка программ с помощью Objective CAML	 61
 VI Приложения	 63

Глава 1

Как получить Objective CAML

Часть I

Основы языка

Глава 2

Функциональное программирование

Глава 3

Императивное программирование

Глава 4

Функциональный и императивный стиль

Глава 5

Графический интерфейс

Глава 6

Приложения

Часть II

Средства разработки

Глава 7

Компиляция и переносимость

Глава 8

Библиотеки

Глава 9

Автоматический сборщик мусора

Глава 10

Средства анализа программ

Глава 11

Средства лексического и синтаксического анализа

11.1 Введение

Определение и реализация средств лексического и синтаксического анализа являлись важным доменом исследования в информатике. Эта работа привела к созданию генераторов лексического и синтаксического анализа `lex` и `yacc`. Команды `camllex` `camlyacc`, которые мы представим в этой главе, являются их достойными наследниками. Два указанных инструмента стали *de-facto* стандартными, однако существуют другие средства, как например потоки или регулярные выражения из библиотеки `Str`, которые могут быть достаточны для простых случаев, там где не нужен мощный анализ.

Необходимость подобных инструментов особенно чувствовалась в таких доменах как компиляция языков программирования. Однако и другие программы могут с успехом использовать данные средства: базы данных, позволяющие определять запросы или электронная таблица, где содержимое ячейки можно определить как результат какой-нибудь формулы. Проще говоря, любая программа, в которой взаимодействие с пользователем осуществляется при помощи языка, использует лексический и синтаксический анализ.

Возьмём простой случай. Текстовый формат часто используется для хранения данных, будь то конфигурационный системный файл или данные табличного файла. Здесь, для использования данных, необходим лексический и синтаксический анализ.

Обобщая, скажем что лексический и синтаксический анализ преобразует линейный поток символов в данные с более богатой структурой:

последовательность слов, структура записи, абстрактное синтаксическое дерево программы и т.д.

У каждого языка есть словарный состав (лексика) и грамматика, которая описывает каким образом эти составные объединяются (синтаксис). Для того, чтобы машина или программа могли корректно обрабатывать язык, этот язык должен иметь точные лексические и синтаксические правила. У машины нет «тонкого чувства» для того чтобы правильно оценить двусмысленность натуральных языков. По этой причине к машине необходимо обращаться в соответствии с чёткими правилами, в которых нет исключений. В соответствии с этим, понятия лексики и семантики получили формальные определения, которые будут кратко представлены в данной главе.

11.2 План главы

Данная глава знакомит нас со средствами лексического и синтаксического анализа, которые входят в дистрибутив Objective CAML. Обычно синтаксический анализ следует за лексическим. В первой части мы узнаем о простом инструменте лексического анализа из модуля `Genlex`. После этого ознакомимся с формализмом рациональных выражений и тем самым рассмотрим более детально определение множества лексических единиц. А так же проиллюстрируем их реализацию в модуле `Str` и инструменте `ocamllex`. Во второй части мы определим грамматику и рассмотрим правила создания фраз языка. После этого рассмотрим два анализа фраз: восходящий и нисходящий. Они будут проиллюстрированы использованием `Stream` и `ocamlyacc`. В приведённых примерах используется контекстно-независимая грамматика. Здесь мы узнаем как реализовать контекстный анализ при помощи `Stream`. В третьей части мы вернёмся к интерпретатору BASIC (см. стр ??) и при помощи `ocamllex` и `ocamlyacc` добавим лексические и синтаксические функции анализа языка.

11.3 Лексика

Синтаксический анализ это предварительный и необходимый этап обработки последовательностей символов: он разделяет этот поток в последовательность слов, так называемых лексические единицы или лексемы.

11.3.1 Модуль Genlex

В данном модуле имеется элементарное средство для анализа символьного потока. Для этого используются несколько категорий предопределённых лексических единиц. Эти категории различаются по типу:

```

1 # type token = Kwd of string
2   | Ident of string
3   | Int of int
4   | Float of float
5   | String of string
6   | Char of char ;;

```

Таким образом мы можем распознать в потоке символов целое число (конструктор `Int`) и получить его значение (аргумент конструктора `int`). Распознаваемые символы и строки подчиняются следующим общепринятым соглашениям: строка окружена символами ("), а символ окружён ('). Десятичное число представлено либо используя запись с точкой (например 0.01), либо с мантиссой и экспонентой (на пример 1E-2). Кроме этого остались конструкторы `Kwd` и `Ident`.

Конструктор `Ident` предназначен для определения идентификаторов. Идентификатором может быть имя переменной или функции языка программирования. Они состоят из любой последовательности букв и цифр, могут включать символ подчёркивания (`_`) или апостроф (`'`). Данная последовательность не должна начинаться с цифры. Любая последовательность следующих операндов тоже будет считаться идентификатором: `+`, `*`, `>` или `-`. И наконец, конструктор `Kwd` определяет категорию специальных идентификаторов или символов.

Категория ключевых слова — единственная из этого множества, которую можно сконфигурировать. Для того, чтобы создать лексический анализатор, воспользуемся следующей конструкцией, которой необходимо передать список ключевых слов на место первого аргумента.

```

1 # Genlex.make_lexer ;;
2 - : string list -> char Stream.t -> Genlex.token Stream.t = <fun>

```

Тем самым получаем функцию, которая принимает на вход поток символов и возвращает поток лексических единиц (с типом `token`).

Таким образом, мы без труда реализуем лексический анализатор для интерпретатора BASIC. Объявим множество ключевых слов:

```

1 # let keywords =
2   [ "REM"; "GOTO"; "LET"; "PRINT"; "INPUT"; "IF"; "THEN"; "—
    ";

```

```

3      "!"; "+" ; "-" ; "*" ; "/" ; "%";
4      "=" ; "<" ; ">" ; "<=" ; ">=" ; "<>";
5      "&" ; "|" ] ;;

```

При помощи данного множества, определим функцию лексического анализа:

```

1 # let line_lexer l = Genlex.make_lexer keywords (Stream.of_string l) ;;
2 val line_lexer : string -> Genlex.token Stream.t = <fun>
3 # line_lexer "LET_x_=_x_+_y_*_3" ;;
4 - : Genlex.token Stream.t = <abstr>

```

Приведённая функция `line_lexer`, из входящего потока символов создаёт поток соответствующих лексем.

11.3.2 Использование потоков

Мы также можем реализовать лексический анализ «в ручную» используя потоки.

В следующем примере определён лексический анализатор арифметических выражений. Функции `lexer` передаётся поток символов из которого она создаёт поток лексических единиц с типом `lexeme Stream.t`¹. Символы пробел, табуляция и переход на новую строку удаляются. Для упрощения, мы не будем обрабатывать переменные и отрицательные целые числа.

```

1 # let rec spaces s =
2     match s with parser
3         [ <" " ; rest > ] -> spaces rest
4         | [ <"\t" ; rest > ] -> spaces rest
5         | [ <"\n" ; rest > ] -> spaces rest
6         | [ <> ] -> ();;
7 val spaces : char Stream.t -> unit = <fun>
8 # let rec lexer s =
9     spaces s;
10    match s with parser
11        [ <"(' ' > ] -> [ <'Lsymbol "(" ; lexer s > ]
12        | [ <"')' > ] -> [ <'Lsymbol ")" ; lexer s > ]
13        | [ <"'+ ' > ] -> [ <'Lsymbol "+" ; lexer s > ]
14        | [ <"'- ' > ] -> [ <'Lsymbol "-" ; lexer s > ]
15        | [ <"'* ' > ] -> [ <'Lsymbol "*" ; lexer s > ]
16        | [ <"'/' > ] -> [ <'Lsymbol "/" ; lexer s > ]

```

¹тип `lexeme` определён на стр. ??


```

17 | [< '0'..'9' as c;
18   i,v = lexint (Char.code c - Char.code('0')) >]
19   -> [<'Lint i ; lexer v>]
20 and lexint r s =
21   match s with parser
22     [< '0'..'9' as c >]
23     -> let u = (Char.code c) - (Char.code '0') in lexint (10*r + u) s
24     | [<>] -> r,s ;;
25 val lexer : char Stream.t -> lexeme Stream.t = <fun>
26 val lexint : int -> char Stream.t -> int * char Stream.t = <fun>

```

Функция `lexint` предназначена для анализа той части потока символов, которая соответствует числовой постоянной. Она вызывается, когда функция `lexer` встречает цифры. В этом случае функция `lexint` поглощает все последовательные цифры и выдаёт соответствующее значение полученного числа.

11.4 Регулярные выражения

Оставим ненадолго практику и рассмотрим проблему лексических единиц с теоретической точки зрения.

Лексическая единица является словом. Слово образуется при конкатенации элементов алфавита. В нашем случае алфавитом является множество символов ASCII.

Теоретически, слово может вообще не содержать символов (пустое слово ²) или состоять из одного символа.

Теоретические исследования конкатенации элементов алфавита для образования лексических элементов (лексем) привели к созданию простого формализма, известного как регулярные выражения.

Определение

Регулярные выражения позволяют определить множества слов. Пример такого множества: идентификаторы. Принцип определения основан на некоторых теоретико-множественных операциях. Пусть M и N — два множества слов, тогда мы можем определить:

- объединение M и N , записываемое $M|N$.
- дополнение, записываемое $\wedge M$: множество всех слов, кроме тех, которые входят в M .

²традиционно, такое слово обозначается греческой буквой эpsilon: ε

- конкатенация M и N : множество всех слов созданных конкатенацией слова из M и слова из N . Записывается просто MN .
- мы можем повторить операцию конкатенации слов множества M и тем самым получить множество слов образованных из конечной последовательности слов множеств M . Такое множество записывается M^+ . Он содержит все слова множества M , все слова полученные конкатенацией двух слов множества M , трёх слов, и т.д. Если мы желаем чтобы данное множество содержало пустое слово, необходимо писать M^* .
- для удобства, существует дополнительная конструкция $M^?$, которая включает все слова множества M , а так же пустое слово.

Один единственный символ ассоциируется с одно элементным множеством. В таком случае выражение $a/b/c$ описывает множество состоящее из трёх слов a , b и c . Существует более компактная запись: $[abc]$. Так как наш алфавит является упорядоченным (по порядку кодов ASCII), можно определить интервал. Например, множество цифр запишется как $[0 - 9]$. Для группировки выражений можно использовать скобки.

Для того, чтобы использовать в записи сами символы-операторы, как обычные символы, необходимо ставить перед ними обратную косую черту: \backslash . Например множество $(\backslash^*)^*$ обозначает множество последовательностей звёздочек.

Пример

Пусть существует множество из цифр $(0,1,2,3,4,5,6,7,8,9)$, символы плюс $(+)$ и минус $(-)$, точки $(.)$ и буквы E . Теперь мы можем определить множество чисел num . Назовем $integers$ множество определённое выражением $[0 - 9]^+$. Множество неотрицательных чисел $unum$ определяется так:

$$integers?(.integers)?(E(\backslash + | -)?integers)?$$

Множество отрицательных и положительных чисел записывается:

$$unum| - unum \quad \text{или} \quad -?unum$$

Распознавание

Теперь, после того как множество выражений определено, остаётся проблема распознавания принадлежности строки символов или одной из её

подстрок этому множеству. Для решения данной задачи необходимо реализовать программу обработки выражений, которая соответствует формальным определениям множества. Для регулярных выражений такая обработка может быть автоматизирована. Подобная автоматизация реализована в модуле **Genlex** из библиотеки **Str** и инструментом **ocamllex**, которые будут представлены в следующих двух параграфах.

Библиотека **Str**

В данном модуле имеется абстрактный тип **regex** и функция **regex**. Указанный тип представляет регулярные выражения, а функция трансформирует регулярное выражение, представленное в виде строки символов, в абстрактное представление.

Модуль **Str** содержит несколько функций, которые используют регулярные выражения и манипулируют символьными строками. Синтаксис регулярных выражений библиотеки **Str** приведён в таблице 11.1.

Таблица 11.1: Регулярные выражения

.	любой символ, кроме \
*	ноль или несколько экземпляров предыдущего выражения
+	хотя бы один экземпляр предыдущего выражения
?	ноль или один экземпляр предыдущего выражения
[..]	множество символов
.	любой символ, кроме \
	интервал записывается при помощи - (пример $[0 - 9]$)
	дополнение записывается при помощи \wedge (пример $[\wedge A - Z]$)
\wedge	начало строки (не путать с дополнением \wedge)
\$	конец строки
	вариант
(..)	группировка в одно выражение (можно ссылаться на это выражение)
i	числовая константа i ссылается на i -ый элемент группированного выражения
\	забой, используется для сопоставления зарезервированных символов в регулярных выражениях

Пример

В следующем примере напомним функцию, которая переводит дату из английского формата во французский. Предполагается, что входной файл

состоит из строк, разбитых на поля данных и элементы даты разделяются точкой. Определим функцию, которая из полученной строки (строка файла), выделяет дату, разбивает её на части, переводит во французских формат и тем самым заменяет старую дату на новую.

```

1  # let french_date_of d =
2    match d with
3      [mm; dd; yy] -> dd ^ "/" ^ mm ^ "/" ^ yy
4      | _ -> failwith "Bad_date_format" ;;
5  val french_date_of : string list -> string = <fun>
6
7  # let english_date_format = Str.regexp "[0-9]+\.[0-9]+\.[0-9]+" ;;
8  val english_date_format : Str.regexp = <abstr>
9
10 # let trans_date l =
11   try
12     let i=Str.search_forward english_date_format l 0 in
13     let d1 = Str.matched_string l in
14     let d2 = french_date_of (Str.split (Str.regexp "\.") d1) in
15     Str.global_replace english_date_format d2 l
16   with Not_found -> l ;;
17   val trans_date : string -> string = <fun>
18
19 # trans_date
20 ".....06.13.99....." ;;
21 - : string = ".....13/06/99....."

```

11.4.1 Инструмент Ocamllex

`ocamllex` — это лексический генератор созданный для Objective CAML по модели `lex`, написанном на языке C. При помощи файла, описывающего элементы лексики в виде множества регулярных выражений, которые необходимо распознать, он создаёт файл-исходник на Objective CAML. К описанию каждого лексического элемента можно привязать какое-нибудь действие, называемое семантическое действие. В полученном коде используется абстрактный тип `lexbuf` из модуля `Lexing`. В данном модуле также имеется несколько функций управления лексическими буферами, которые могут быть использованы программистом для того чтобы определить необходимые действия.

Обычно, файлы, описывающие лексику, имеют расширение `.mll`. Для того, чтобы из файла `lex_file.mll` получить файл на Objective CAML,

необходимо выполнить следующую команду:

```
1 ocamllex lex_file.mll
```

После этого, мы получим файл `lex_file.ml`, содержащий код лексического анализатора. Теперь, данный файл можно использовать в программе на Objective CAML. Каждому множеству правил анализа соответствует функция, которая принимает лексический буфер (типа `Lexing.lexbuf`) и затем возвращает значение, определённое семантическим действием. Значит, все действия для определённого правила должны создавать значение одного и того же типа.

Формат у файла для `ocamllex` следующий:

```
1 {
2   header
3 }
4 let ident = regexp
5     ...
6 rule ruleset1 = parse
7     regexp { action }
8     | ...
9     | regexp { action }
10 and ruleset2 = parse
11     ...
12 and ...
13 {
14   trailer-and-end
15 }
```

Части «заголовок» и «продолжение-и-конец» не являются обязательными. Здесь вставляется код Objective CAML, определяющий типы данных, функции и т.д. необходимые для обработки данных. В последней части используются функции, которые используют правила анализа множества лексических данных из средней части. Серия объявлений, которая предшествует определению правил, позволяет дать имя некоторым регулярным выражениям. Эти имена будут использоваться в определении правил.

Пример

Вернёмся к нашему интерпретатору BASIC и усовершенствуем тип возвращаемых лексических единиц. Таким образом, мы можем воспользоваться функцией `lexer`, (см. стр. ??) которая возвращает такой же тип результата (`lexeme`), но на входе она получает буфер типа `Lexing.lexbuf`.

```

1  {
2    let string_chars s =
3      String.sub s 1 ((String.length s)-2) ;;
4  }
5
6  let op_ar = ['-','+', '*', '\%', '/',]
7  let op_bool = ['!', '&', '|']
8  let rel = ['=', '<', '>']
9
10 rule lexer = parse
11   [' '] { lexer lexbuf }
12   | op_ar { Lsymbol (Lexing.lexeme lexbuf) }
13   | op_bool { Lsymbol (Lexing.lexeme lexbuf) }
14   | "<=" { Lsymbol (Lexing.lexeme lexbuf) }
15   | ">=" { Lsymbol (Lexing.lexeme lexbuf) }
16   | "<>" { Lsymbol (Lexing.lexeme lexbuf) }
17   | rel { Lsymbol (Lexing.lexeme lexbuf) }
18   | "REM" { Lsymbol (Lexing.lexeme lexbuf) }
19   | "LET" { Lsymbol (Lexing.lexeme lexbuf) }
20   | "PRINT" { Lsymbol (Lexing.lexeme lexbuf) }
21   | "INPUT" { Lsymbol (Lexing.lexeme lexbuf) }
22   | "IF" { Lsymbol (Lexing.lexeme lexbuf) }
23   | "THEN" { Lsymbol (Lexing.lexeme lexbuf) }
24   | '-'? ['0'-'9']+ { Lint (int_of_string (Lexing.lexeme lexbuf)) }
25   | ['A'-'z']+ { Lident (Lexing.lexeme lexbuf) }
26   | "'[_^_']"* '[_^_'] { Lstring_(string_chars_(Lexing.lexeme_lexbuf))_ }

```

После обработки данного файла командой `ocamllex` получим функцию `lexer` типа `Lexing.lexbuf -> lexeme`. Далее, мы рассмотрим каким образом подобные функции используются в синтаксическом анализе (см. стр. ??).

Глава 12

Взаимодействие с языком С

Глава 13

Приложения

Часть III

Устройство программы

Глава 14

Модульное программирование

Глава 15

Объектно-ориентированное программирование

Глава 16

Сравнение моделей устройств программ

Глава 17

Приложения

Часть IV

Параллелизм и распределение

Глава 18

Процессы и связь между процессами

18.1 Параллельное программирование

Глава 19

Распределённое программирование

Глава 20

Приложения

Часть V

Разработка программ с помощью Objective CAML

Часть VI

Приложения

