

# Optimizing Pattern Matching

Fabrice Le Fessant, Luc Maranget

INRIA Roquencourt, B.P. 105, 78153 Le Chesnay Cedex, France  
(Email: {*Fabrice.Le\_fessant, Luc.Maranget*}@inria.fr)

## ABSTRACT

We present improvements to the backtracking technique of pattern-matching compilation. Several optimizations are introduced, such as commutation of patterns, use of exhaustiveness information, and control flow optimization through the use of labeled static exceptions and context information. These optimizations have been integrated in the Objective-Caml compiler. They have shown good results in increasing the speed of pattern-matching intensive programs, without increasing final code size.

## 1. INTRODUCTION

Pattern-matching is a key feature of functional languages. It allows to discriminate between the values of a deeply structured type, binding subparts of the value to variables at the same time. ML users now routinely rely on their compiler for such a task; they write complicated, nested, patterns. And indeed, transforming high-level pattern-matching into elementary tests is a compiler job. Moreover, because it considers the matching as a whole and that it knows some intimate details of runtime issues such as the representation of values, compiler code is often better than human code, both as regards compactness and efficiency.

There are two approaches to pattern-matching compilation, the underlying model being either decision trees [5] or backtracking automata [1]. Using decision trees, one produces *a priori* faster code (because each position in a term is tested at most once), while using backtracking automata, one produces *a priori* less code (because patterns never get copied, hence never get compiled more than once). The price paid in each case is losing the advantage given by the other technique.

This paper mostly focuses on producing faster code in the backtracking framework. Examining the code generated by the Objective-Caml compiler [11], which basically used the Augustsson's original algorithm, on small frequently found programs, such as a list-merge function, or on large examples [14], we found that the backtracking scheme could still

be improved.

Our optimizations improve the produced backtracking automaton by grouping elementary tests more often, removing useless tests and avoiding the blind backtracking behavior of previous schemes. To do so, the compiler uses new information and outputs a new construct. New information include incompatibility between patterns, exhaustiveness information and contextual information at the time of backtracking. As to the new construct, previous schemes used a lone “exit” construct whose effect is to jump to the nearest enclosing “trap-handler” ; we enrich both exits and trap-handlers with labels, resulting in finer control of execution flow.

Our optimizations also apply to or-patterns, a convenient feature to group clauses with identical actions. Unsharing of actions is avoided by using our labelled exit construct. As or-patterns may contain variables, the exit construct is also extended to take arguments.

All our optimizations are now implemented in the latest version of the Objective-Caml compiler, whose language of accepted patterns has been extended by allowing variables in or-patterns.

The structure of this article is the following: we first introduce some theoretical basics on pattern-matching in section 2 and describe the compilation scheme to backtracking automata in section 3. Then, we briefly introduce our optimizations and or-pattern compilation in an intuitive way in sections 4 and 5, while section 6 is a formalization of our complete compilation scheme. Finally, some experimental results are shown in section 7, and a comparison with other approaches is discussed in section 8.

## 2. BASICS

In this section, we introduce some notations and definitions. Most of the material here is folklore, save, perhaps, or-patterns.

### 2.1 Patterns and Values

ML is a typed language, where new types of values can be introduced using *type definitions* such as:

```
type t = Nil | One of int | Cons of int * t
```

This definition introduces a type `t`, with three *constructors* that build values of type `t`. These three constructors define the *complete signature* of type `t`. Every constructor has an arity, i.e. the number of arguments it takes. Here arity of `Nil` is zero, while the arities of `One` and `Cons` are one and two respectively. A constructor of arity zero is called

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICFP'01, September 3-5, 2001, Florence, Italy.

Copyright 2001 ACM 1-58113-415-0/01/0009 ...\$5.00.

a *constant* constructor, while other constructors are *non-constant* constructors.

Most native data types in ML – such as integers, records, arrays, tuples – can be seen as particular instances of such type definitions. For example, in the following we will consider lists (*nil* being the constant constructor  $[]$  and *cons* the infix constructor  $::$ ), and tuples (the type of  $n$ -tuples defines one constructor of arity  $n$ , pairs being written with the infix constructor “,”). For our purpose, integers are constant constructors, and the signature of the integer type is infinite.

More formally, *patterns* and *values* are defined as follows:

$p ::=$	<b>Patterns</b>
$\_$	wildcard
$x$	variable
$c(p_1, p_2, \dots, p_a)$	constructor pattern
$(p_1   p_2)$	or-pattern
$v ::=$	<b>Values</b>
$c(v_1, v_2, \dots, v_a)$	constructor value

In the following, we freely replace variables by wild-cards “ $\_$ ” when their names are irrelevant. While describing compilation, convenient tools are vectors of values ( $\vec{v} = (v_1 v_2 \dots v_n)$ ) and  $\vec{v}_{n \leftrightarrow m} = (v_n \dots v_m)$ , vectors of patterns ( $\vec{p} = (p_1 p_2 \dots p_n)$ ) and  $\vec{p}_{n \leftrightarrow m} = (p_n \dots p_m)$ ) and matrices of patterns ( $P = (p_j^i)$ ).

In this paper, we present pattern-matching compilation as a transformation on an intermediate code in the compiler, called *lambda-code*. Here, another useful object is the *clause matrix* ( $P \rightarrow L$ ):

$$(P \rightarrow L) = \left( \begin{array}{cccc} p_1^1 & p_2^1 & \dots & p_n^1 & \rightarrow & l^1 \\ p_1^2 & p_2^2 & \dots & p_n^2 & \rightarrow & l^2 \\ & & & \vdots & & \\ p_1^m & p_2^m & \dots & p_n^m & \rightarrow & l^m \end{array} \right)$$

A clause matrix associates rows of patterns ( $p_1^i p_2^i \dots p_n^i$ ) to lambda-code actions  $l^i$ .

## 2.2 Pattern Matching in ML

A pattern can be seen as representing a set of values sharing a common prefix.

**DEFINITION 1 (INSTANCE).** *Let  $p$  be a pattern and  $v$  be a value belonging to a common type. The value  $v$  is an instance of the pattern  $p$  or  $p$  matches  $v$ , written  $p \preceq v$  when one of the following rules apply:*

$$\begin{array}{ll} \_ & \preceq v \\ x & \preceq v \\ (p_1 | p_2) & \preceq v \quad \text{iff } p_1 \preceq v \text{ or } p_2 \preceq v \\ c(p_1, \dots, p_a) & \preceq c(v_1, \dots, v_a) \quad \text{iff } (p_1 \dots p_a) \preceq (v_1 \dots v_a) \\ (p_1 \dots p_a) & \preceq (v_1 \dots v_a) \quad \text{iff } p_i \preceq v_i, \forall i \in [1..a] \end{array}$$

Seeing a pattern as the set of its instances, it is clear that or-patterns express set union.

In ML, patterns are a binding construct, more specifically, a successful match  $p \preceq v$ , binds the variables of  $p$  to some sub-terms of  $v$ . Such bindings can be computed while checking that  $p$  matches  $v$ , provided that the following set  $\mathcal{V}(p)$

of variables defined by  $p$  is well-defined:

$$\begin{aligned} \mathcal{V}(\_) &= \emptyset \\ \mathcal{V}(x) &= \{x\} \\ \mathcal{V}(c(p_1, \dots, p_a)) &= \mathcal{V}(p_1 \dots p_a) \\ \mathcal{V}(p_1 \dots p_a) &= \mathcal{V}(p_1) \cup \dots \cup \mathcal{V}(p_a) \\ &\quad \text{if for all } i \neq j, \mathcal{V}(p_i) \cap \mathcal{V}(p_j) = \emptyset \\ \mathcal{V}(p_1 | p_2) &= \mathcal{V}(p_1), \text{ if } \mathcal{V}(p_1) = \mathcal{V}(p_2) \end{aligned}$$

The first “if” condition above is the linearity of patterns. The second condition is specific to or-patterns, it means that matching by either side of the or-pattern binds the same variables (additionally, homonymous variables should possess the same type).

We then define the now dominant, textual priority scheme to disambiguate the case when several rows in a matrix match:

**DEFINITION 2 (MATCHING PREDICATE).** *Let  $P$  be a pattern matrix and  $\vec{v} = (v_1 \dots v_n)$  be a value vector. The value  $v$  matches line number  $i$  in  $P$ , if and only if the following two conditions are satisfied:*

- $(p_1^i \dots p_n^i) \preceq (v_1 \dots v_n)$
- $\forall j < i, (p_1^j \dots p_n^j) \not\preceq (v_1 \dots v_n)$

We will not give a full semantics for evaluating pattern-matching expressions, and more generally lambda-code. Intuitively, given a clause matrix  $P \rightarrow L$  and a value vector  $\vec{v}$  such that line number  $i$  in  $P$  matches  $\vec{v}$ , evaluating the matching of  $\vec{v}$  by  $P \rightarrow L$  in some environment  $\rho$  is evaluating  $l^i$  in  $\rho$  extended by the bindings introduced while matching  $\vec{v}$  by  $(p_1^i \dots p_n^i)$ . If  $\vec{v}$  is not matched by any line in  $P$ , we say that the pattern-matching  $P$  *fails*. If no such  $\vec{v}$  exists, then the pattern-matching is said *exhaustive*.

Like pattern vectors, pattern matrices represent sets of value vectors. More specifically, when some line in  $P$  matches  $\vec{v}$  we simply say that  $P$  matches  $\vec{v}$ . This looks obvious, but representing sets using matrices is at the core of our optimizations. One easily checks that the instances of  $P$  are the union of the instances of the lines of  $P$ . That is, when considering a matrix globally, the second condition in definition 2 above is irrelevant. More important, row order is also irrelevant.

Finally, the instance relation induces relations on the patterns themselves.

**DEFINITION 3 (RELATIONS ON PATTERNS).** *We define the following three relations:*

1. Pattern  $p$  is less precise than pattern  $q$ , written  $p \preceq q$ , when all instances of  $q$  are instances of  $p$ .
2. Pattern  $p$  and  $q$  are equivalent, written  $p \equiv q$ , when their instances are the same.
3. Patterns  $p$  and  $q$  are compatible when  $p$  and  $q$  share a common instance.

Here some remarks are to be made. Because of typing, checking the precision relation is neither obvious nor cheap. More precisely, there is no simple way to decide whether  $p \preceq \_$  holds or not. For instance,  $([] | \_ :: \_) \preceq \_$  holds, while  $(\text{Nil} | \text{One } \_) \preceq \_$  does not. Or-patterns are not responsible for this complication, since we also have  $(\_, \_) \preceq \_$ . In such cases one should “expand”  $p$  and consider, whether

signatures are complete or not (see [13, Section 5.1]). By contrast, compatibility can be checked by a simple recursive algorithm. When compatible, patterns  $p$  and  $q$  admit a least upper bound, written  $p \uparrow q$ , which can be computed while checking compatibility:

$$\left\{ \begin{array}{l} (p_1 \dots p_a) \uparrow (q_1 \dots q_a) = (p_1 \uparrow q_1 \dots p_a \uparrow q_a) \\ \quad \quad \quad \uparrow q = q \\ \quad \quad \quad p \uparrow - = p \\ c(p_1, \dots, p_a) \uparrow c(q_1, \dots, q_a) = c(r_1, \dots, r_a) \\ \quad \quad \quad \text{where } (r_1 \dots r_a) \text{ is } (p_1 \dots p_a) \uparrow (q_1 \dots q_a) \end{array} \right.$$

With the following additional rules for or-patterns:

$$(p_1 | p_2) \uparrow q = \begin{cases} p_1 \uparrow q, & \text{when } p_2 \text{ and } q \text{ not compatible} \\ p_2 \uparrow q, & \text{when } p_1 \text{ and } q \text{ not compatible} \\ (p_1 \uparrow q | p_2 \uparrow q), & \text{otherwise} \end{cases}$$

$$p \uparrow (q_1 | q_2) = (q_1 | q_2) \uparrow p$$

Proving that  $p \uparrow q$  is indeed the least upper bound of  $p$  and  $q$  is easy, by considering patterns as sets of their instances. Note that  $p \uparrow q$  is defined up to  $\equiv$ -equivalence, and that it encodes instance intersection.

### 3. COMPILATION

In this section, we present a compilation scheme close to the one described in [20, 1], and implemented in compilers such as the `hbc` compiler or the Objective Caml compiler. This classical scheme will be refined later into an optimized scheme, using same notations and concepts.

#### 3.1 Output of the match compiler

The compilation of pattern-matching is described by the scheme  $\mathcal{C}$  that maps a clause matrix to a *lambda-code* expression. We now describe the specific *lambda-code* constructs that the scheme  $\mathcal{C}$  outputs while compiling patterns.

- Let-bindings: `let` ( $x$   $l_x$ )  $l$ , nested let-bindings are abbreviated as:

`let` ( $x_1$   $l_1$ ) ( $x_2$   $l_2$ )  $\dots$  ( $x_n$   $l_n$ )  $l$

- Static exceptions, `exit` and traps, `catch`  $l_1$  `with`  $l_2$ . If, when evaluating the body  $l_1$ , `exit` is encountered, then the result of evaluating `catch`  $l_1$  `with`  $l_2$  is the result of evaluating the handler  $l_2$ , otherwise it is the result of evaluating  $l_1$ . By contrast with dynamic exceptions, static exceptions are directly compiled as jumps to the associated handlers (plus some environment adjustment, such as stack pops), whereas traps do not generate any code.

- Switch constructs:

`switch`  $l$  `with`  
`case`  $c_1$ :  $l_1 \dots$  `case`  $c_k$ :  $l_k$   
`default`:  $d$

The result of a switch construct is the evaluation of the  $l_i$  corresponding to the constructor  $c_i$  appearing as the head of the value  $v$  of  $l$ . If the head constructor of  $v$  doesn't appear in the case list, the result is the evaluation of the default  $d$  expression.

The default clause `default`:  $d$  can be omitted. In such a case the switch behavior is unspecified on non-recognized values. Scheme  $\mathcal{C}$  can thus omit the default clause when

it is known that case lists will cover all possibilities at runtime. We use the keyword `switch*` to highlight switch constructs with no default clause.

Those switch constructs are quite sophisticated, they compile later into more basic constructs: tests, branches and jump tables. We in fact modified the Objective Caml compiler to improve the compilation of switch constructs, using techniques first introduced in the context of compiling the case statement of Pascal [3]. The key points are using range tests, which can typically be performed by one single (unsigned) test and branch plus possibly one addition, cutting sparse case lists into denser ones, and deciding which of jump tables or test sequence is more appropriate to each situation. A survey of these techniques can be found in [19].

- Accessors: `field`  $n$   $x$ , where  $x$  is a variable and  $n$  is an integer offset. By convention, the first argument of non-constant constructors stands at offset zero.
- Sequences:  $l_1$ ;  $l_2$  and units:  $()$

#### 3.2 Initial state

Input to the pattern matching compiler  $\mathcal{C}$  consists of two arguments: a vector of variables  $\vec{x}$  of size  $n$  and a clause matrix  $P \rightarrow L$  of width  $n$  and height  $m$ .

$$\vec{x} = (x_1 \ x_2 \ \dots \ x_n), \quad P \rightarrow L = \begin{pmatrix} p_1^1 & p_2^1 & \dots & p_n^1 & \rightarrow & l^1 \\ p_1^2 & p_2^2 & \dots & p_n^2 & \rightarrow & l^2 \\ & & & \vdots & & \\ p_1^m & p_2^m & \dots & p_n^m & \rightarrow & l^m \end{pmatrix}$$

The initial matrix is generated from source input. Given a pattern-matching expression (in Caml syntax):

`match`  $x$  `with` |  $p^1 \rightarrow e^1$  |  $p^2 \rightarrow e^2 \dots$  |  $p^m \rightarrow e^m$

The initial call to  $\mathcal{C}$  is:

`catch`  
 $\mathcal{C}((x), \begin{pmatrix} p^1 & \rightarrow & l^1 \\ p^2 & \rightarrow & l^2 \\ & & \rightarrow \\ p^m & \rightarrow & l^m \end{pmatrix})$   
`with` (`failwith` "Partial match")

Where the  $l^i$ 's are the translations to lambda-code of the  $e^i$ 's, and (`failwith` "Partial match") is a runtime failure that occurs when the whole pattern matching fails.

#### 3.3 Classical scheme

By contrast with previous presentations, we assume that matrix  $P \rightarrow L$  has at least one row (i.e.  $m > 0$ ). This condition simplifies our presentation, without restricting its generality. Hence, scheme  $\mathcal{C}$  is defined by cases on non-empty clause matrices:

1. If  $n$  is zero (i.e. when there are no more columns), then the first row of  $P$  matches the empty vector  $()$ :

$$\mathcal{C}(), \begin{pmatrix} \rightarrow & l^1 \\ \rightarrow & l^2 \\ & \vdots \\ \rightarrow & l^m \end{pmatrix} = l^1$$

2. If  $n$  is not zero, then a simple compilation is possible, using the following four rules.

- (a) If all patterns in the first column of  $p$  are variables,  $y^1, y^2, \dots, y^m$ , then:

$$\mathcal{C}(\vec{x}, P \rightarrow L) = \mathcal{C}((x_2 \ x_3 \ \dots \ x_n), P' \rightarrow L')$$

where

$$P' \rightarrow L' = \left( \begin{array}{ccc} p_2^1 & \cdots & p_n^1 \rightarrow \text{let } (y^1 \ x_1) \ l^1 \\ p_2^2 & \cdots & p_n^2 \rightarrow \text{let } (y^2 \ x_1) \ l^2 \\ & \vdots & \\ p_2^m & \cdots & p_n^m \rightarrow \text{let } (y^m \ x_1) \ l^m \end{array} \right)$$

We call this rule, the *variable* rule. This case also handles wild-card patterns: they are treated like variables except that the let-binding is omitted.

- (b) If all patterns in the first column of  $P$  are constructor patterns  $c(q_1, \dots, q_a)$ , then let  $C$  be the set of *matched constructors*, that is, the set of the head constructors of the  $p_1^i$ 's.

Then, for each constructor  $c$  in  $C$ , we define the *specialized* clause matrix  $\mathcal{S}(c, P \rightarrow L)$  by mapping the following transformation on the rows of  $P$ .

$p_1^i$	$\mathcal{S}(c, P \rightarrow L)$
$c(q_1^i, \dots, q_a^i)$	$q_1^i \cdots q_a^i \ p_2^i \cdots p_n^i \rightarrow l_i$
$c'(q_1^i, \dots, q_a^i)$ ( $c' \neq c$ )	No row

(Matrices  $\mathcal{S}(c, P \rightarrow L)$  and  $P \rightarrow L$  define the same matching predicate when  $x_1$  is bound to some value  $c(v_1, \dots, v_a)$ .) Furthermore, for a given constructor  $c$  of arity  $a$ , let  $y_1, \dots, y_a$  be fresh variables. Then, for any constructor  $c$  in  $C$ , we define the lambda-expression  $r(c)$ :

$$\begin{aligned} & (\text{let } (y_1 \ (\text{field } 0 \ x_1)) \\ & \quad \dots \\ & \quad (y_a \ (\text{field } (a-1) \ x_1)) \\ & \quad \mathcal{C}((y_1, \dots, y_a, x_2, \dots, x_n), \mathcal{S}(c, P \rightarrow L))) \end{aligned}$$

Finally, assuming  $C = \{c_1, \dots, c_k\}$ , the compilation result is:

$$\begin{aligned} & \text{switch } x_1 \ \text{with} \\ & \quad \text{case } c_1 : r(c_1) \cdots \text{case } c_k : r(c_k) \\ & \quad \text{default: exit} \end{aligned}$$

(Note that the default clause can be omitted when  $C$  makes up a full signature.) We call this rule, the *constructor* rule.

- (c) If  $P$  has only one row and that this row starts with an or-pattern:

$$P = ((q_1 \mid \dots \mid q_o) \ p_2 \ \cdots \ p_n \rightarrow l),$$

Then, compilation result is:

$$\mathcal{C}((x_1), \left( \begin{array}{c} q_1 \rightarrow () \\ \vdots \\ q_o \rightarrow () \end{array} \right)); \mathcal{C}((x_2 \ \dots \ x_n), (p_2 \ \dots \ p_n \rightarrow l))$$

This rule is the *orpat* rule. Observe that it does not duplicate any pattern nor action. However, variables in

or-patterns are not supported, since, in clause  $q_i \rightarrow ()$ , the scope of  $q_i$  variables is the action “()”.

- (d) Finally, if none of the previous rules applies, the clause matrix  $P \rightarrow L$  is cut in two clause matrices  $P_1 \rightarrow L_1$  and  $P_2 \rightarrow L_2$ , such that  $P_1 \rightarrow L_1$  is the largest prefix of  $P \rightarrow L$  for which one of the variable, constructor or orpat rule applies.

Then, compilation result is:

$$\text{catch } \mathcal{C}(\vec{x}, P_1 \rightarrow L_1) \ \text{with } \mathcal{C}(\vec{x}, P_2 \rightarrow L_2)$$

This rule is the *mixture* rule.

This paper doesn't deal with optimizing let-bindings, which are carelessly introduced by scheme  $\mathcal{C}$ . This job is left to a later compilation phase.

## 4. OPTIMIZATIONS

We now describe some improvement to the classical compilation scheme. For simplicity, we present examples and defer the full presentation of our scheme to section 6. In all these examples, we focus on pattern-matching compilation, replacing potentially arbitrary actions more simple ones, such as integers or variables.

### 4.1 Optimizing the mixture rule

In this section and in the following, our running example is the classical list-merge:

```
let merge lx ly = match lx,ly with
| [], _ -> 1
| _, [] -> 2
| x::xs, y::ys -> 3
```

Such a matching on pairs encodes matching on two arguments. As a consequence, we consider the following initial call to scheme  $\mathcal{C}$ :

$$\mathcal{C}((\text{lx ly}), (P \rightarrow L))$$

Where  $(P \rightarrow L)$  is:

$$(P \rightarrow L) = \left( \begin{array}{cc} [] & \_ \rightarrow 1 \\ \_ & [] \rightarrow 2 \\ x::xs & y::ys \rightarrow 3 \end{array} \right)$$

Applying the mixture rule twice yields three matrices:

$$\begin{aligned} P_1 \rightarrow L_1 &= \left( \begin{array}{cc} [] & \_ \rightarrow 1 \\ \_ & [] \rightarrow 2 \end{array} \right) \\ P_2 \rightarrow L_2 &= \left( \begin{array}{cc} \_ & [] \rightarrow 2 \\ x::xs & y::ys \rightarrow 3 \end{array} \right) \\ P_3 \rightarrow L_3 &= \left( \begin{array}{cc} x::xs & y::ys \rightarrow 3 \end{array} \right) \end{aligned}$$

Now, consider another clause matrix  $(P' \rightarrow L')$ :

$$(P' \rightarrow L') = \left( \begin{array}{cc} [] & \_ \rightarrow 1 \\ x::xs & y::ys \rightarrow 3 \\ \_ & [] \rightarrow 2 \end{array} \right)$$

Both clause matrices define the same matching function, namely they both map  $([] \ v)$  to 1,  $(v_1::v_2 \ [])$  to 2 and  $(v_1::v_2 \ v'_1::v'_2)$  to 3. Furthermore,  $(P' \rightarrow L')$  can be obtained from  $(P \rightarrow L)$  by swapping its second and third row. More generally, one easily checks that swapping two contiguous *incompatible* rows is legal. Then applying the mixture rule to  $(P' \rightarrow L')$ , yields two matrices only:

$$\begin{aligned} P'_1 \rightarrow L'_1 &= \left( \begin{array}{cc} [] & \_ \rightarrow 1 \\ x::xs & y::ys \rightarrow 3 \end{array} \right), \\ P'_2 \rightarrow L'_2 &= \left( \begin{array}{cc} \_ & [] \rightarrow 2 \end{array} \right) \end{aligned}$$

```

catch
  (catch
    (switch lx with case []: 1
      default: exit)
    with (catch
      (switch ly with case []: 2
        default: exit)
      with (catch
        (switch lx with
          case (::):
            (switch ly with
              case (::) : 3
                default: exit)
            default: exit))))
    with (failwith "Partial match"))

```

```

catch
  (catch* lx with
    case []: 1
    case (::) :
      (switch ly with
        case (::): 3
        default: exit))
  with
    (switch ly with
      case []: 2
      default: exit)
  with (failwith "Partial match")

```

Figure 1: Mixture optimization

Final outputs for  $P \rightarrow L$  and  $P' \rightarrow L'$  are displayed on Figure 1. Hence, as a result of replacing  $P \rightarrow L$  by  $P' \rightarrow L'$ , the two tests on `lx` that were performed separately on the left code are now merged in a single switch in the right code. Also notice that one trap disappears.

More generally, an optimized mixture rule should take advantage of pattern-matching semantics to swap rows when possible, so that as few cuts as possible are performed.

## 4.2 Using exhaustiveness information

The Objective Caml compiler checks the exhaustiveness of pattern matching expressions and issues a warning before compiling non-exhaustive pattern matchings. However, the exhaustiveness information can also be used for avoiding tests. Matrix  $P'$  of the previous section is exhaustive; this means that there will be no "Partial match" failure at runtime. As an immediate consequence, the switch: `(switch ly with case []: 2 default: exit)` always succeeds (this switch is the last one performed by the optimized code in figure 1). Thus, we replace it by 2. We can also suppress the outermost trap. Hence, applying both optimizations described up to now, compilation of  $P \rightarrow L$  finally yields:

```

catch
  (switch* lx with
    case []: 1
    case (::): (switch ly with
              case []: 3
              default: exit))
  with 2

```

In the general case, exhaustiveness information is exploited by slightly modifying scheme  $\mathcal{C}$ . It suffices to avoid emitting default clauses in switch constructs, when it is known that no exit should escape from produced code. This property holds initially for exhaustive pattern matchings, and transmits to all recursive calls, except for the call on  $P_1 \rightarrow L_1$  in the mixture rule.

## 4.3 Optimizing exits

The two previous optimizations yield optimal code for the merge example. Hence we complicate the running example by considering a matching on objects of type `t` from sec-

tion 2:

$$P \rightarrow L = \begin{pmatrix} \text{Nil} & - & \rightarrow 1 \\ - & \text{Nil} & \rightarrow 2 \\ \text{One } x & - & \rightarrow 3 \\ - & \text{One } y & \rightarrow 4 \\ \text{Cons } (x, xs) & \text{Cons } (y, ys) & \rightarrow 5 \end{pmatrix}$$

The optimized mixture rule yields four matrices:

$$\begin{aligned} P_1 \rightarrow L_1 &= \begin{pmatrix} \text{Nil} & - & \rightarrow 1 \\ \text{Cons } (x, xs) & \text{Cons } (y, ys) & \rightarrow 5 \end{pmatrix} \\ P_2 \rightarrow L_2 &= \begin{pmatrix} - & \text{Nil} & \rightarrow 2 \end{pmatrix} \\ P_3 \rightarrow L_3 &= \begin{pmatrix} \text{One } x & - & \rightarrow 3 \end{pmatrix} \\ P_4 \rightarrow L_4 &= \begin{pmatrix} - & \text{One } y & \rightarrow 4 \end{pmatrix} \end{aligned}$$

For reasons that will appear immediately, we apply the mixture rule from bottom to top, thereby nesting trap handlers. The match being exhaustive, compilation yields the code displayed on the left part of Figure 2.

Now, consider what happens at run-time when `(lx ly)` is `(Cons (v1, v2) One v)`. A first switch on `lx` leads to line 7, where a switch on `ly` is performed. This switch fails, and the default action jumps to the nearest enclosing handler (line 13), where `ly` is tested against `Nil` resulting in another switch failure. Here, in our case, control goes to line 17, where another switch on `lx` (against `One x`) fails, resulting in final jump to line 20.

Hence, it would be appropriate to jump to line 20 right from the first test on `ys`. To do so, both exits and trap handlers are now *labelled* by integers. Note that this new feature does not really complicate the compilation of static exceptions. Then, it becomes possible to jump to different trap handlers from the same point and a better compilation of  $P \rightarrow L$  is displayed in the right part of figure 2.

The code above maps vectors `(Cons (v1, v2) One v)` to 4 by executing two switches, while previous code needed four switches to perform the same task. Hence, exit optimization has a noticeable benefit as regards run-time efficiency. As regards code size, exit optimization may increase it, since some switches may have larger case lists. However, code size remains under control, since no extra switches are generated. Hence, final code size critically depends on how switches translate to machine-level constructs. For instance, machine-level code size obviously does not increase when

```

1 catch
2   (catch
3     (catch
4       (switch lx with
5         case Nil: 1
6         case Cons:
7           (switch ly with
8             case Cons: 5
9             default: exit)
10        default: exit)
11      with
12        (switch ly with
13          case Nil: 2
14          default: exit))
15      with
16        (switch lx with
17          case One: 3
18          default: exit))
19    with 4
20  with 4

```

Unoptimized code

```

1 catch
2   (catch
3     (catch
4       (switch lx with
5         case Nil: 1
6         case Cons:
7           (switch* ly with
8             case Cons: 5
9             case Nil: (exit 2)
10            case One: (exit 4))
11          default: (exit 2))
12      with (2)
13        (switch ly with
14          case Nil: 2
15          default: (exit 3)))
16      with (3)
17        (switch lx with
18          case One: 3
19          default: (exit 4))
20    with (4) 4

```

Optimized code

Figure 2: Exit optimization

switches are translated to jump tables<sup>1</sup>.

Surprisingly, performing exit optimization is quite simple and cheap: the needed information is available at compile-time by inspecting pattern matrices only. Reachable trap handlers are defined as pairs  $(P, e)$  of a pattern matrix and an integer. Reachable trap handlers originate from the division performed by the mixture rule. Here,  $P_1 \rightarrow L_1$  is compiled with the reachable trap-handlers  $(P_2, 2)$ ,  $(P_3, 3)$  and  $(P_4, 4)$ . Then, the constructor rule specializes reachable trap handlers. Here, in the case where  $lx$  is  $Cons(v_1, v_2)$ , specializing reachable trap handlers results in  $((Nil), 2)$  and  $((One\ y), 4)$  (note that specializing  $P_3$  yields an empty matrix, which is discarded). Hence, while generating the first switch on  $ly$  (line 7), it is known that the code produced by compiling trap handlers number 2 and 3 will surely exit when  $ly$  is  $One\ v$ , and a jump to trap handler number 4 can be generated by the compiler in that case.

#### 4.4 Aggressive control flow optimization

The code produced by exit optimization still contains redundant tests, some of which can be removed without altering the handler structure introduced by the mixture rule. More specifically, we consider trap handler number 3 (line 16). It results from compiling  $P_3$  and is a switch of  $lx$  against  $One$ .

The only `(exit 3)` lies in trap handler number 2 (line 15) and results from  $ly$  not being  $Nil$ , this gives us no direct information on  $lx$ . Now, looking upwards for `(exit 2)`, we can infer that trap handler number 2 is entered from two different points. In the first case (line 9),  $(lx\ ly)$  is fully known as  $(Cons(v_1, v_2)\ Nil)$ , in the second case (line 11), only  $lx$  is known to be  $One\ v$ . As `(exit 3)` on line 15 gets executed only when  $ly$  is not  $Nil$ , we can finally deduce that

<sup>1</sup>Given the Objective Caml encoding of constructors, we are here in the same desirable situation where the compilation of apparently larger switches does not result in producing more code.

the first case never results in entering trap handler number 3. As a consequence, trap handler number 3 is executed in a context where  $lx$  necessarily is  $One\ v$ , the switch it performs is useless and line 16 can be simplified into “3”. This elimination of useless tests[4] is usually performed at a lower level by combining dead code elimination[9] and conditional constant propagation[21, 6].

Finally, after all optimizations, there remains one redundant switch in produced code, in trap-handler number 2 (line 12). As a result, vectors  $(Cons(v_1, v_2)\ Nil)$  are mapped to 2 by testing  $ly$  twice. One should notice that this is precisely the test that would get duplicated by compilation to decision trees.

Describing what is known on values while entering trap handlers is slightly involved. The key idea is representing set of value vectors as pattern matrices. We call such a set a *context*. Contexts for the three trap handlers of our example are:

Trap number	Context
2	$\begin{pmatrix} One\ \_ & \_ \\ Cons(\_, \_) & Nil \end{pmatrix}$
3	$\begin{pmatrix} One\ \_ & (One\ \_   Cons(\_, \_)) \end{pmatrix}$
4	$\begin{pmatrix} Cons(\_, \_) & One\ \_ \end{pmatrix}$

If precise enough and exploited fully, we conjecture that contexts subsume exhaustiveness information. However as intuition suggests and experience confirms, contexts get larger while compilation progresses, potentially reaching huge sizes at the end of matrices. We cure this by safely approximating contexts when they get too large, replacing some patterns in them by wild-cards. Hence the optimizations of section 4.2 is still worth considering, as being cheap and always applicable.

## 5. COMPILING OR-PATTERNS

Until now, the code produced for or-patterns is inefficient, because only one or-pattern can be compiled at a time, re-

quiring multiple applications of the mixture rule before and after each or-pattern. Thanks to integer labelled exits, one easily avoids dividing matrices before or-patterns. Consider a “*car*” function for our three-constructors list:

```
let car list = match list with
| Nil -> -1
| (One x | Cons (x, _)) -> x
```

Compilation proceeds by allocating a new trap-handler number 2 and expanding the clause “One x | Cons (x,\_)” into two clauses with patterns “One x” and “Cons (x,\_)”. Actions for the new clauses are exits to 2:

```
catch
  C((list), (
    Nil      -> -1
    One x'   -> (exit 2 x')
    Cons (x', _) -> (exit 2 x')
  ))
with (2 x) C(( ), ( -> x ))
```

Note that both exits and trap handlers now take yet another extra argument, the occurrences of  $x'$  in exits are non-binding and refer to pattern variables, while the occurrence of  $x$  in handler is binding. This new construct allows the compilation of or-patterns with variables. Implementation is not very tricky: the `catch ... with (2 x) ...` construct allocates one mutable variable; an exit updates this variable, which is read before entering the handler. In a native code compiler, such a variable is a temporary and ultimately a machine register. The generated lambda-code is as follow:

```
catch
  switch* list with
  case Nil: -1
  case One: (exit 2 (field 0 list))
  case Cons: (exit 2 (field 0 list))
with (2 x) x
```

Moreover, by the semantics of pattern-matching, cuts after or-patterns can also be avoided in many situations. In the case of one column matrices, where the expanded or-patterns express the full matching performed, all cuts can be avoided. Things get a bit more complicated when matrices have more than one column. Consider the following clause matrix,

$$P \rightarrow L = \begin{pmatrix} (1|2) & p_2 \rightarrow l^1 \\ (3|4) & q_2 \rightarrow l^2 \end{pmatrix}$$

We further assume a match on  $(x\ y)$  and that match failure should result in (exit 1) (the static exception label corresponding to match failure can be given as a third argument to the compilation scheme). Writing  $p_1 = (1|2)$  and  $q_1 = (3|4)$ , there are obviously no value vectors  $(v_1\ v_2)$  such that  $v_1$  is an instance of both  $p_1$  and  $q_1$ . As a consequence, the following compilation is correct:

```
catch
  (catch
    (switch x with
      case 1: (exit 2) case 2: (exit 2)
      case 3: (exit 3) case 4: (exit 3)
      default: (exit 1))
    with (2) C((y), ( p_2 -> l^1 ), 1))
  with (3) C((y), ( q_2 -> l^2 ), 1)
```

Intuitively, once  $x$  is checked, the choice between first and second row is made. Depending on the value of  $y$ , matching may still fail, but then, the whole matching fails.

Conversely, matrix division cannot be avoided when matching by  $p_1$  does not exclude matching by  $q_1$ , that is, when  $p_1$  and  $q_1$  are compatible. This is the case, for instance, when  $p_1 = (1|2)$  and  $q_1 = (2|3)$ . Then, a correct compilation is:

```
catch
  (catch
    (switch x with
      case 1: (exit 2) case 2: (exit 2)
      default: (exit 3))
    with (2) C((y), ( p_2 -> l^1 ), 3))
  with (3)
    (catch
      (switch x with
        case 2: (exit 4) case 3: (exit 4)
        default: (exit 1))
      with (4) C((y), ( q_2 -> l^2 ), 1))
```

Note that the third argument to the first recursive call to the compilation scheme is “3” and not “1”. As a consequence, vectors  $(2\ v_2)$  such that  $p_2$  does not match  $v_2$  while  $q_2$  matches  $v_2$  get mapped correctly to  $l^2$ . A slight inefficiency shows up, since  $x$  is tested twice. More striking, perhaps, vectors  $(1\ v_2)$  such that  $p_2$  does not match  $v_2$  also lead to testing  $x$  twice.

An alternative compilation rule for or-pattern would simply expand or-patterns in a pre-processing phase, yielding the matrix:

$$\begin{pmatrix} 1 & p_2 \rightarrow l^1 \\ 2 & p_2 \rightarrow l^1 \\ 2 & q_2 \rightarrow l^2 \\ 3 & q_2 \rightarrow l^2 \end{pmatrix}$$

Then, there are no extra run-time tests on  $x$ , since the constructor rule applies. However, patterns  $p_2$  and  $q_2$  are now compiled twice. Note that there is no simple solution for avoiding this duplication of effort, since, once the constructor rule is applied, the two occurrences of these patterns occur in different contexts. More generally, code size is now out of control, a clear contradiction with the spirit of backtracking automata.

## 6. OUR COMPILATION SCHEME

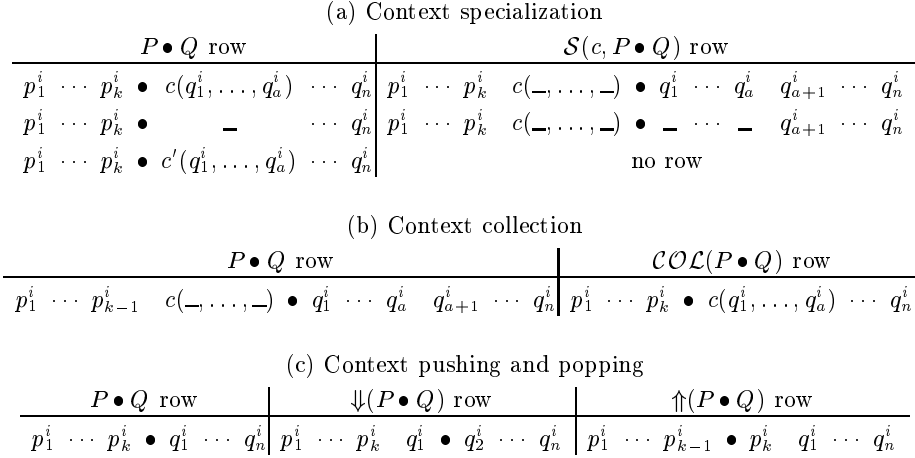
The new scheme  $\mathcal{C}^*$  takes five arguments and a typical call is  $\mathcal{C}^*(\vec{x}, P \rightarrow L, ex, def, ctx)$ , where  $\vec{x} = (x_1 \dots x_n)$  and  $P \rightarrow L$  is a clause matrix of width  $n$ :

$$P \rightarrow L = \begin{pmatrix} p_1^1 & \dots & p_n^1 & \rightarrow & l^1 \\ p_1^2 & \dots & p_n^2 & \rightarrow & l^2 \\ & & & & \vdots \\ p_1^m & \dots & p_n^m & \rightarrow & l^m \end{pmatrix}$$

Extra arguments are:

- The exhaustiveness argument  $ex$  is either *partial* or *total* depending on whether compilation can produce escaping exit constructs or not.
- Reachable trap handlers  $def$  are sequences  $(P_1, e_1); \dots; (P_t, e_t)$ , where the  $e_i$ 's are integers (trap handler numbers) and the  $P_i$ 's are pattern matrices of width  $n$ .

**Figure 3: Operations on contexts**



- The context  $ctx$  is a pattern matrix of width  $k+n$ , equivalent to a pair of matrixes  $P \bullet Q$ , where each row is divided into a prefix (in  $P$ ) of width  $k$  and a fringe (in  $Q$ ) of width  $n$ .

$$P \bullet Q = \begin{pmatrix} p_1^1 & \cdots & p_k^1 & \bullet & q_1^1 & \cdots & q_n^1 \\ p_1^2 & \cdots & p_k^2 & \bullet & q_1^2 & \cdots & q_n^2 \\ \vdots & & \vdots & & \vdots & & \vdots \\ p_1^m & \cdots & p_k^m & \bullet & q_1^m & \cdots & q_n^m \end{pmatrix}$$

Informally, at any point in compilation, contexts are pre-order representations of what is known about matched values. The fringe records the possible values for  $\vec{x}$ , while the prefix records the same information for other sub-terms which are relevant to pending calls to  $\mathcal{C}^*$ . Transfers of patterns from fringe to prefix are performed on the arguments of recursive calls, while transfers in the opposite direction are performed as results are collected.

The initial call to  $\mathcal{C}^*$  for an exhaustive match is:

$$\mathcal{C}^*((x), \begin{pmatrix} p^1 & \rightarrow & l^1 \\ p^2 & \rightarrow & l^2 \\ \vdots & & \vdots \\ p^m & \rightarrow & l^m \end{pmatrix}, total, \emptyset, (\bullet \_))$$

For a non-exhaustive match,  $ex$  is *partial*,  $def$  is the one-element sequence  $((\_, 1))$  and a trap handler is added as in section 3.3. The context argument remains the same: it expresses that nothing is known yet about the value of  $\vec{x}$ .

The new scheme returns a lambda-code  $l$  and a *jump summary*,  $\rho = \{\dots, i \mapsto ctx, \dots\}$ , which is a mapping from trap numbers to contexts. Jump summaries describe what is known about matched values at the places where (`exit i ...`) occur in  $l$ .

## 6.1 Operations on contexts

We define the following four operations on contexts :

- (a) Context *specialization*,  $\mathcal{S}$ , by a constructor  $c$  of arity  $a$  is defined by mapping the transformation of figure 3-(a) on context rows.

- (b) Context *collection*,  $\mathcal{COL}$ , is the reverse of specialization. It combines the the last element of the prefix with the appropriate number of arguments standing at beginning of the fringe (see figure 3-(b)).

- (c) Context *pushing*  $\Downarrow$  and *popping*  $\Uparrow$  move the fringe limit one step forward and backward, without examining any pattern (see figure 3-(c)).

As contexts are used to represent set of values, we naturally define union and intersection over contexts. Context *union*  $P \bullet Q \cup P' \bullet Q'$  yields a new matrix whose rows are the rows of  $P \bullet Q$  and  $P' \bullet Q'$ . Row order is not relevant. Context intersection  $P \bullet Q \cap P' \bullet Q'$  is defined as a context whose rows are the least upper bounds of the compatible rows of  $P \bullet Q$  and  $P' \bullet Q'$ . Context *extraction*  $\mathcal{EX}$  is a particular case of context intersection.

$$\mathcal{EX}(p, P' \bullet Q') = (\_ \dots \_ \bullet p \dots \_) \cap P' \bullet Q'$$

For example, when  $p$  is  $c(\_, \dots, \_)$ , context extraction retains those value vectors represented by  $P' \bullet Q'$  whose  $k+1^{\text{th}}$  components admit  $c$  as head constructor. Observe that such a computation involves extracting or-pattern arguments and making wild-cards more precise.

Except for collection and popping, which consume prefix elements, all these operations can be extended to simple matrices, by using an empty prefix in input, and taking the fringe for output. Doing so, we obtain exactly the operations of section 3.3 used to compute pattern matrices (specialization  $\mathcal{S}$  in particular).

Operations on contexts are extended to jump summaries in the natural manner. For instance, the union of  $\rho$  and  $\rho'$  is defined as:

$$\rho \cup \rho' = \{\dots, i \mapsto \rho(i) \cup \rho'(i), \dots\}$$

Operations on matrices are extended to reachable trap handlers in a similar manner: for instance, pushing trap handlers is defined as pushing all matrices in them :

$$\Downarrow((P_1, e_1); \dots; (P_t, e_t)) = (\Downarrow(P_1), e_1); \dots; (\Downarrow(P_t), e_t)$$



## 6.2 Compilation scheme

We now describe scheme  $\mathcal{C}^*$  by considering cases over the typical call.

1. If  $n$  is zero, then we have:

$$\mathcal{C}^*((), \left( \begin{array}{c} \rightarrow l^1 \\ \rightarrow l^2 \\ \vdots \\ \rightarrow l^m \end{array} \right), ex, def, ctx) = l^1, \emptyset$$

Observe that the jump summary is empty since no exit is outputed.

2. With respect to section 3.3, the variable rule only changes as regards the extra arguments  $ex$ ,  $def$  and  $ctx$ . We only describe these changes. The performed recursive call returns code  $l$  and jump summary  $\rho$  :

$$l, \rho = \mathcal{C}^*(\dots, \dots, ex, \Downarrow(def), \Downarrow(ctx))$$

Exhaustiveness information  $ex$  does not change, while  $def$  and  $ctx$  are pushed.

The variable rule returns  $l$  unchanged and  $\rho$  popped.

3. In the constructor rule, let  $C = \{c_1, \dots, c_k\}$  be the matched constructors, let also  $\Sigma$  be the signature of their type. For a given constructor  $c \in C$ , the performed recursive call is:

$$\mathcal{C}^*(\dots, \dots, ex, \mathcal{S}(c, def), \mathcal{S}(c, ctx))$$

Exhaustiveness information  $ex$  is passed unchanged, while the other two extra arguments are specialized (specialization of trap handlers being the natural extension of matrix specialization).

Each recursive call returns a lambda-code  $l(c)$  and a jump summary  $\rho_c$ . Lambda-code  $l(c)$  gets wrapped into let-bindings like in section 3.3, yielding the final lambda-code  $r(c)$ . We then define a case list  $\mathcal{L}$  and a jump summary  $\rho_{\text{rec}}$  as follows:

$$\begin{aligned} \mathcal{L} &= \text{case } c_1 : r(c_1) \cdots \text{case } c_k : r(c_k) \\ \rho_{\text{rec}} &= \{ \dots, i \mapsto \bigcup_{c \in C} \mathcal{C}\mathcal{O}\mathcal{L}(\rho_c(i)), \dots \} \end{aligned}$$

The case list is as before, while the jump summary is the union of the the jump summaries produced by recursive calls, once collected.

Optimizations are then performed. For clarity, optimizations are described as a two phase process: first, extend (or not extend) the case list  $\mathcal{L}$  with constructors taken from  $\Sigma \setminus C$ , and add (or not add) a default case; then, compute the final jump summary.

A first easy case is when  $\Sigma \setminus C$  is empty or when  $ex$  is *total*. Then, the case list  $\mathcal{L}$  is not augmented. Otherwise, we distinguish two cases :

- (a) If  $\Sigma \setminus C$  is finite, then for all constructors  $c$  in this set we consider the context

$$Q_c \bullet Q'_c = \mathcal{E}\mathcal{X}(c(\_, \dots, \_), ctx)$$

Then, trap handlers  $(P_1, e_1); \dots; (P_t, e_t)$  are scanned left-to-right, stopping at the smallest  $i$ , such that the intersection  $Q'_c \cap P_i$  is not empty. That is, we find the

trap handler where to jump to when the head constructor of  $x_1$  is  $c$ , in order to extend the case list as follows :

$$\mathcal{L} = \mathcal{L} \text{ case } c : (\text{exit } e_i)$$

It is possible that  $e_i$  does not exist (when  $Q'_c$  is empty). This means that  $x_1$  head constructor will never be  $c$  at runtime.

- (b) If  $\Sigma \setminus C$  is infinite (as in the case of integers) or considered too large (as it might be in the case of characters), then, a default case is added to the case list :

$$\mathcal{L} = \mathcal{L} \text{ default: } (\text{exit } e_1)$$

That is, all non-recognized constructors lead to a jump the nearest enclosing reachable trap-handler.

However it is still possible to extend the case list for particular constructors, applying the previous procedure (a) to the constructors that appear in the first column of reachable trap handler matrices and not in  $C$ .

The final jump summary is computed by considering the final case list  $\mathcal{L}$ . For a given trap handler number  $e_i$  let  $\{c'_1, \dots, c'_{k'}\}$  be the set of constructors such that **case**  $c'_j$  : (**exit**  $e_i$ ) appears in  $\mathcal{L}$ . Then the jump summary  $\rho_{e_i}$  is defined as:

$$\rho_{e_i} = \{ e_i \mapsto \mathcal{E}\mathcal{X}(c'_1(\_, \dots, \_) \mid \cdots \mid c'_{k'}(\_, \dots, \_)), ctx \}$$

Moreover, if there is a default clause, the jump summary  $\rho_d$  is defined as:

$$\rho_d = \{ e_1 \mapsto ctx \}$$

Finally the constructor rule returns a switch on case list  $\mathcal{L}$  and the jump summary built by performing the union of  $\rho_{\text{rec}}$ , of all  $\rho_{e_i}$ 's and, when appropriate, of  $\rho_d$ .

The constructor rule performs many context unions, so that contexts may become huge. Fortunately, contexts can be made smaller using a simple observation. Namely, let  $\vec{p}$  and  $\vec{q}$  be two rows in a context, such that  $\vec{p}$  is less precise than  $\vec{q}$  (i.e., all instances of  $\vec{q}$  are instances of  $\vec{p}$ ). Then, row  $\vec{q}$  can be removed from the context, without modifying its meaning as a set of value vectors. Hence, while performing context union, one can leave aside some pattern rows. If the produced context is still too large, then contexts are safely approximated by first replacing some patterns in them by wild-cards (typically all the pattern in a given column) and then removing rows using the previous remark. Rough experiments lead us to set the maximal admissible context size to 32 rows, yielding satisfactory compilation time in pathological examples and exact contexts in practical examples.

4. Or-pattern compilation operates on matrices whose first column contains at least one or-pattern. Additionally, when  $p_1^i$  is a or-pattern, then for all  $j$ ,  $i < j \leq m$  one of the following, mutually exclusive, conditions must hold:

- (a)  $p_1^i$  and  $p_1^j$  are *not* compatible.
- (b)  $p_1^i$  and  $p_1^j$  are compatible, and  $(p_2^i \dots p_n^i)$  is less precise than  $(p_2^j \dots p_n^j)$

Conditions (a) and (b) guarantee that, whenever  $p_1^i$  matches the first value vector  $v_1$  of a value  $\vec{v}$ , but row  $i$  does not match  $\vec{v}$ , then no further row in  $P$  matches  $\vec{v}$  either. This

is necessary since further rows of  $P$  won't be reachable in case of failure in the or-pattern trap handler.

Now, consider one row number  $i$ , such that  $p_1^i$  is the or-pattern  $q_1 \mid \dots \mid q_o$ . Further assume that this or-pattern binds the variables  $y_1, \dots, y_v$ . First, we allocate a fresh trap number  $e$  and divide  $P \rightarrow L$  into the following or-body  $P' \rightarrow L'$  and or-trap  $P'' \rightarrow L''$  clauses:

$$P' \rightarrow L' = \left( \begin{array}{cccc} \vdots & & & \\ p_1^{i-1} & \dots & p_n^{i-1} & \rightarrow & l^{i-1} \\ q_1 & \dots & - & \rightarrow & (\text{exit } e \ y_1 \dots y_v) \\ \vdots & & & & \\ q_o & \dots & - & \rightarrow & (\text{exit } e \ y_1 \dots y_v) \\ p_1^{i+1} & \dots & p_n^{i+1} & \rightarrow & l^{i+1} \\ \vdots & & & & \end{array} \right)$$

$$P'' \rightarrow L'' = ( p_2^i \ \dots \ p_m^i \rightarrow l^i )$$

In the or-body matrix, observe that the or-pattern is expanded, while the other patterns in row number  $i$  are replaced by wild-cards and the action is replaced by exits.

Recursive calls are performed as follows:

$$\begin{aligned} l', \rho' &= \mathcal{C}^*(\vec{x}, P' \rightarrow L', ex, def, ctx) \\ l'', \rho'' &= \dots \\ \dots \mathcal{C}^*(\vec{x}_{2 \leftrightarrow n}, P'' \rightarrow L'', ex, \downarrow(\mathcal{E}\mathcal{X}(p, def)), \downarrow(\mathcal{E}\mathcal{X}(p, ctx))) \end{aligned}$$

Outputted code finally is `catch l' with (e y1... yv) l''` and the returned jump summary is  $\rho = \rho' \cup \uparrow(\rho'')$ .

- The mixture rule is responsible for feeding the other rules with appropriate clause matrices. We first consider the case of a random division. Hence let us cut  $P \rightarrow L$  into  $Q \rightarrow M$  and  $R \rightarrow N$  at some row. Then a fresh trap number  $e$  is allocated and a first recursive call is performed:

$$l_q, \rho_q = \mathcal{C}^*(\vec{x}, Q \rightarrow M, \text{partial}, (R, e); def, ctx)$$

The exhaustiveness information is *partial*, since nothing about the exhaustiveness of  $Q$  derives from the exhaustiveness of  $P$ . Reachable trap handlers are extended.

Then, a second recursive call is performed:

$$l_r, \rho_r = \mathcal{C}^*(\vec{x}, R \rightarrow N, ex, def, \rho_q(e))$$

It is no surprise that the context argument to the new call is extracted from the jump summary of the previous call. Argument  $ex$  does not change. Indeed, if matching by  $P$  cannot fail, then matching by  $R$  neither can.

Then, the scheme can output the code

$$l = \text{catch } l_q \text{ with } (e) \ l_r$$

and return the jump summary  $(\rho_q \setminus \{e\}) \cup \rho_r$ , where  $\rho_q \setminus \{e\}$  stands for  $\rho_q$  with the binding for  $e$  removed.

Of course, our optimizing compiler does not perform a random division into two matrices. It instead divides  $P \rightarrow L$  right away into several sub-matrices. This can be described formally as several, clever, applications of the random mixture rule, so that one of the three previous rules apply to each matrix in the division. The aim of the optimizing mixture rule is thus to perform a division of  $P$  into as few sub-matrices as possible. We present a simple, greedy, approach that scans  $P$  downwards.

We only describe the case when  $p_1^1$  is a constructor pattern. Thus, having performed the classical mixture rule, we are in a situation where the  $i$  topmost rows of  $P$  have a constructor pattern in first position (i.e. are constructor rows for short) and where  $p_1^{i+1}$  is not a constructor pattern. At that point, a matrix  $C$  has been built, which encompasses all the rows of  $P$  from 1 to  $i$ . Let us further write  $P'$  for what remains of  $P$ , and let  $O$  and  $R$  be two new, initially empty matrices. We then scan the rows of  $P'$  from top to bottom, appending them at the end of  $C$ ,  $O$  or  $R$ . That is, given row number  $j$  in  $P'$ :

- If  $p_1^j$  is a variable, then append row  $j$  at the end of  $R$ .
- If  $p_1^j$  is a constructor pattern, then ...
  - If row  $j$  is not compatible with all the rows of both  $R$  and  $O$ , then append row  $j$  at the end of  $C$  (i.e., move row  $j$  above all the rows that have been extracted from  $P'$  at previous stages).
  - If row  $j$  is not compatible with all the rows of  $R$  and that one of conditions (a) or (b) for applying the or-pattern rule are met by  $O$  with row  $j$  appended at the end, then do such an append.
  - Otherwise, append row  $j$  at the end of  $R$ .
- If  $p_1^j$  is an or-pattern, then consider cases (ii) and (iii).

When the scan of  $P'$  is over, three matrices,  $C$ ,  $O$  and  $R$  have been built. In the case where  $O$  is empty, matrix  $C$  is valid input to the constructor rule; otherwise, appending the rows of  $O$  at the end of  $C$  yields valid input for applying (maybe more than once) the or-pattern rule, which will in turn yield valid input to the constructor rule (provided that  $(\_ \mid \dots)$  or patterns have been replaced by semantically equivalent wild-cards in a previous phase). Thus, the matrix built by appending  $O$  at the end of  $C$  is recorded into the overall division and the division process is restarted with input  $R$ , unless  $R$  is empty.

Finally, the full process divides the input matrix  $P$  into several matrices, each of which is valid input to the other rules of the compilation scheme.

## 7. EXPERIMENTAL RESULTS

We compare the performance of the code generated by the Objective-Caml compilers version 3.00 and 3.01, where the former implements the scheme of section 3.3 and the latter implements our new optimizing scheme (there are other differences of minor relevance to our purpose). For most programs there is little difference; this is natural since pattern-matching usually accounts for a small fraction of most programs running time. A full analysis of the efficiency of our optimizations would in fact require counting relevant instructions (test, branches and indirect branches through jump tables), both statically and dynamically. By lack of time, we only present some programs that demonstrate significant improvement.

Our first benchmark is the traditional `fib`, that we write using an or-pattern.

```
let rec fib n = match n with
| (0|1) -> 1 | _ -> fib (n-1) + fib (n-2)
```

Here, we simply measure the execution time of computing `fib 38`. Our second benchmark, `pcf`, is a byte-code compiler and interpreter for PCF. We compute the geometric

mean of the execution time for a set of five different PCF programs. The time-consuming part of this program is the byte-code machine which we coded in the style of the byte-code machine included in [14], the winning entry of the 2000 ICFP programming contest. (we also give figures for this program under the name `raytrace`).

Experiments were performed on a lightly loaded 366Mhz Pentium Pro Linux PC. The tables show wall-clock times (in seconds) and ratios:

	fib		raytrace		pcf	
V 3.00	5.36	100	1.69	100	8.12	100
V 3.01	3.74	71	1.62	96	5.08	63

Obviously, as demonstrated by the `fib` example, compilation of or-patterns has much improved. Testing similar examples confirms that fact. Improvements also comes from the better compilation of switches. The `pcf` example is more interesting, it shows that our optimizations yield a 37% speed-up, in the case of a typical ML application (a quickly written, compact, prototype implementation of some programming language). The `raytrace` example exhibits less important improvements on the whole test suite of the contest; however, improvements are noticeable for some inputs.

It should also be noticed that the new compiler somehow equates the runtime performance of various coding styles, a feature that is important for a high-level construct such as pattern-matching. Variations in coding style include the relative ordering of non-overlapping patterns and on the order of arguments in pairs.

We also performed measurements on a 500Mhz Dec Alpha server. They suggest that the effects of our optimization do not depend on the targeted architecture.

	fib		pcf	
V 3.00	3.4	100	4.13	100
V 3.01	2.5	74	2.86	69

The `raytrace` example is omitted because it relies on IEEE floating point arithmetic, which is not implemented in the Objective Caml compiler for this architecture.

More detailed information on these benchmarks is available at <http://caml.inria.fr/pattern/speed.html>.

## 8. RELATED WORK

### 8.1 Decision Trees vs Backtracking

Compiling to decision trees is the original approach to pattern matching compilation; it first appeared in the Hope compiler and is described in [5]. It is currently used in the SML-NJ compiler [7].

In this approach, there is no mixture rule: instead, the *constructor* rule applies as soon as there is at least one constructor in the first column, and a specialization matrix is created for each matched constructor, plus one additional matrix for the remaining constructors in the signature of the types of matched values, if any. Specialization is done by following the rules of section 6.1. This means that rows whose first pattern is a variable get copied several times.

On the one hand, this approach guarantees that one constructor test is never performed twice. On the other hand, copied pattern rows are compiled independently and this result in potentially large code size. Namely, examples exist that make the SML-NJ compiler produce exponential code [12].

Compilation to backtracking automata is the classical scheme of section 3.3 (see also [1, 20]). It is currently in use in the Haskell-HBC and Objective-Caml compiler [11]. As we already argued, its main advantage is that patterns are never copied, yielding linear output size. Of course, the price paid is potentially testing the same sub-term several times, resulting in potentially poor runtime performance. In that aspect, our new compilation scheme shows that this price can be reduced significantly.

Compilation to decision trees easily detects unused match cases and non-exhaustive matchings, since there is no dead code in a decision tree. Detecting these situations is important, as programmers should be warned about them. However, those problems are NP-complete [17] and this gives us a hint about the potential size of decision trees. More concretely, a decision tree may have many leaves corresponding to non-matched values, whereas knowing that one such values exist is the needed information. Rather, we check unused match cases and exhaustiveness before compilation with a simple algorithm [13] that solves the used matched case problem by basically traversing the decision tree without generating it. Advantages are not generating the tree, stopping search as soon as used match cases are found and applying various heuristics and matrix simplifications which are not relevant to direct compilation. Then, one of our optimizations uses exhaustiveness information.

### 8.2 Compiling or-patterns

From available ML or Haskell compilers, we only found two compilers dealing with or-patterns: the (old) Objective-Caml compiler and the SML-NJ compiler. Our technique makes the old Objective-Caml scheme (see section 3.3) obsolete, by both producing more efficient code and allowing variables in or-patterns.

The SML-NJ approach is very simple to understand and implement: or-patterns are expanded during a pre-processing phase. However, as we already discussed at the end of section 5, this may lead to many duplications of patterns. Such a risk is compatible with the very philosophy of compilation to decision trees and is natural in that context.

### 8.3 Optimizations

Most optimizations dealing with pattern-matching in the literature try to improve the order in which tests are performed. In the matrix-based description, one considers alternatives to systematically choosing the first column of matrices in the constructor rule. Hence, such an approach can be characterized as “column optimization”, while our approach would rather be “row optimization”. Since choosing the best column is thought to be NP-complete (to our knowledge, there is no published proof), most approaches describe heuristics. A typical and early work on such heuristics is [2], a more recent and thorough study is [16]. Another, related in practice, approach relies on sequentiality theory to identify *directions* that are columns that must be tested by all possible matchers [10, 15, 17, 13]. However, computing directions is expansive, and one can consider relying on cheaper heuristics.

These works rather apply to the decision trees, with a primary focus on reducing code size. It is unclear to us how to combine column and row optimization in practice and whether this would yield noticeable improvements or not.

There also exists a partial-evaluation based approach to

pattern-matching optimization. [8] and later [18] specialize an ultra-naive pattern-matching interpreter to create an efficient pattern-matching compiler. Both authors use context information as we do. By contrast, their target is decision trees. In the end, the automatic process of partial evaluation does not find as many optimizations as we do.

## 9. CONCLUSION

This paper contribution is twofold. First, we propose an improvement on the classical technique of compiling pattern-matching expressions into backtracking automata, a technique that has remained virtually the same for about 15 years. Our improvements yield automata which run faster, thereby alleviating the disadvantage of backtracking automata in practical cases. Moreover the very structure of the produced automata is not altered and hence the highly desirable property that output size is linear in the input size is preserved. As a second contribution, we propose a technique for efficiently compiling or-patterns with variables, still preserving the linearity of output size. Using or-patterns in place of “catch-all” wild-cards results in more robust programs, while using one clause with a or-pattern in place of several clauses with identical actions results in more compact, sometime clearer, programs. ML programmers can now enjoy these benefits, without being afraid of degraded runtime efficiency or code size explosion.

We would have wished to make a clear statement on comparing backtracking automata and decision trees. However, sophisticated compilation techniques exist that minimize the drawbacks of both approaches. Those are our techniques for backtracking automata, and hash-consing and column optimizations for decision trees. In the absence of a practical comparison of full-fledged algorithms, choosing one technique or the other reflects one's commitment to guaranteed code size or guaranteed runtime performance.

## 10. REFERENCES

- [1] AUGUSTSSON, L. Compiling pattern matching. In *Functional Programming Languages and Computer Architecture*, J.-P. Jouannaud, Ed. Springer-Verlag, Berlin, DE, 1985, pp. 368–381. Lecture Notes in Computer Science 201 Proceedings of. Conference at Nancy.
- [2] BAUDINET, M., AND MACQUEEN, D. Tree pattern matching for ML. unpublished paper, Dec. 1985.
- [3] BERNSTEIN, R. L. Producing good code for the case statement. *Software—Practice and Experience* 15, 10 (Oct. 1985), 1021–1024.
- [4] BODÍK, R., GUPTA, R., AND SOFFA, M. L. Interprocedural conditional branch elimination. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI-97)* (New York, June 15–18 1997), vol. 32, 5 of *ACM SIGPLAN Notices*, ACM Press, pp. 146–158.
- [5] CARDELLI, L. Compiling a functional language. In *Conference Record of the 1984 ACM Symposium on Lisp and Functional Programming* (Aug. 1984), ACM, ACM, pp. 208–217.
- [6] FRASER, C. W. A compact, machine-independent peephole optimizer. In *Conference Record of the Sixth Annual ACM Symposium on Principles of Programming Languages* (Jan. 1979), ACM, ACM, pp. 1–6.
- [7] HARPER, R. W., MACQUEEN, D. B., AND MILNER, R. Standard ML. Report ECS-LFCS-86-2, Department of Computer Science, University of Edinburgh, Edinburgh, UK, 1986. Also CSR-209-86.
- [8] JØRGENSEN, J. Generating a pattern matching compiler by partial evaluation. In *Glasgow Workshop on Functional Programming, Ullapool* (Glasgow University, July 1990), P. C. J. van Rijsbergen, Ed., Springer-Verlag, pp. 177–195.
- [9] KNOOP, J., RÜTHING, O., AND STEFFEN, B. Partial dead code elimination. In *Proceedings of the Conference on Programming Language Design and Implementation* (New York, NY, USA, June 1994), ACM Press, pp. 147–158.
- [10] LAVILLE, A. Implementation of lazy pattern matching algorithms. In *ESOP'88* (1988), H. Ganzinger, Ed., vol. 300, pp. 298–316.
- [11] LEROY, X. The objective caml system: Documentation and user's manual, 2000. With Damien Doligez, Jacques Garrigue, Didier Rémy, and Jôme Vouillon. Available from <http://caml.inria.fr>.
- [12] MARANGET, L. Compiling lazy pattern matching. In *Proc. of the 1992 conference on Lisp and Functional Programming* (1992), ACM Press.
- [13] MARANGET, L. Two techniques for compiling lazy pattern matching. Research Report 2385, INRIA Rocquencourt, Oct. 1994.
- [14] PLCLUB, AND CAML'R US. Objective-caml: Winner of the first and second prizes of the programming contest. ACM SIGPLAN International Conference on Functional Programming (ICFP '2000).
- [15] PUEL, L., AND SUÁREZ, A. Compiling pattern matching by term decomposition. *Journal of Symbolic Computation* 15, 1 (Jan. 1993), 1–26.
- [16] SCOTT, K., AND RAMSEY, N. When do match-compilation heuristics matter? Tech. Rep. CS-2000-13, Department of Computer Science, University of Virginia, May 2000.
- [17] SEKAR, R. C., RAMESH, R., AND RAMAKRISHNAN, I. V. Adaptive pattern matching. In *Automata, Languages and Programming, 19th International Colloquium* (Vienna, Austria, 13–17 July 1992), W. Kuich, Ed., vol. 623 of *Lecture Notes in Computer Science*, Springer-Verlag, pp. 247–260.
- [18] SESTOFT, P. ML pattern match compilation and partial evaluation, 1996.
- [19] SPULER, D. A. Compiler code generation for multiway branch statements as a static search problem. Tech. Rep. 94/3, Department of Computer Science, James Cook University, 1994.
- [20] WADLER, P. Compilation of pattern matching. In *The Implementation of Functional Programming Languages*, S. L. Peyton Jones, Ed. Prentice-Hall International, 1987, ch. 7.
- [21] WEGMAN, M., AND ZADECK, F. K. Constant propagation with conditional branches. In *Conference Record of the 12th Annual ACM Symposium on Principles of Programming Languages* (New Orleans, LS, Jan. 1985), B. K. Reid, Ed., ACM Press,

pp. 291-299.