

Translation Verification of the pattern matching compiler

Francesco Mecca

%%

1 Introduction

2 Introduction

This dissertation presents an algorithm for the translation validation of the OCaml pattern matching compiler. Given a source program and its compiled version the algorithm checks whether the two are equivalent or produce a counter example in case of a mismatch.

Our equivalence algorithm works with decision trees. Source patterns are converted into a decision tree using a matrix decomposition algorithm. Target programs, described in the Lambda intermediate representation language of the OCaml compiler, are turned into decision trees by applying symbolic execution.

2.1 Translation validation

A pattern matching compiler turns a series of pattern matching clauses into simple control flow structures such as `if`, `switch`, for example:

```
#+BEGIN_SRC ocaml
match x with
| [] -> (0, None)
| x::[] -> (1, Some x)
| _::y::_ -> (2, Some y)
#+END_SRC

(if scrutinee
  (let (field_1 =a (field 1 scrutinee))
```

```

      (if field_1
        (let
          (field_1_1 =a (field 1 field_1)
            x =a (field 0 field_1))
          (makeblock 0 2 (makeblock 0 x)))
        (let (y =a (field 0 scrutinee))
          (makeblock 0 1 (makeblock 0 y))))))
[0: 0 0a])

```

%% TODO: side by side The code in the right is in the Lambda intermediate representation of the OCaml compiler. The Lambda representation of a program is shown by calling the `ocamlc` compiler with `-drawlambda` flag.

The OCaml pattern matching compiler is a critical part of the OCaml compiler in terms of correctness because any bug would result in wrong code production rather than triggering compilation failures. Such bugs also are hard to catch by testing because they arise in corner cases of complex patterns which are typically not in the compiler test suite.

The OCaml core developers group considered evolving the pattern matching compiler, either by using a new algorithm or by incremental refactorings of its codebase. For this reason we want to verify that new implementations of the compiler avoid the introduction of new bugs and that such modifications don't result in a different behaviour than the current one.

One possible approach is to formally verify the pattern matching compiler implementation using a machine checked proof. Another possibility, albeit with a weaker result, is to verify that each source program and target program pair are semantically correct. We chose the latter technique, translation validation because is easier to adopt in the case of a production compiler and to integrate with an existing codebase. The compiler is treated as a blackbox and proof only depends on our equivalence algorithm.

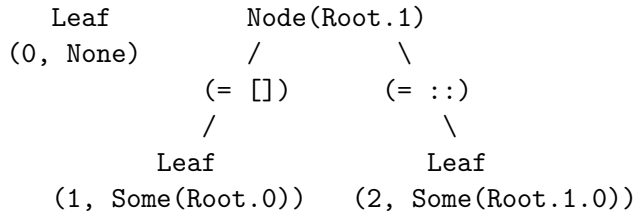
2.2 Our approach

%% replace common TODO Our algorithm translates both source and target programs into a common representation, decision trees. Decision trees were chosen because they model the space of possible values at a given branch of execution. Here is the decision tree for the source example program.

```

      Node(Root)
      /      \
    (= [])   (= ::)
    /        \

```



(Root.0) is called an *accessor*, that represents the access path to a value that can be reached by deconstructing the scrutinee. In this example Root.0 is the first subvalue of the scrutinee.

Target decision trees have a similar shape but the tests on the branches are related to the low level representation of values in Lambda code. For example, cons cells `x:xs` are blocks with tag 0.

To check the equivalence of a source and a target decision tree, we proceed by case analysis. If we have two terminals, such as leaves in the previous example, we check that the two right-hand-sides are equivalent. If we have a node N and another tree T we check equivalence for each child of N , which is a pair of a branch condition π_i and a subtree C_i . For every child (π_i, C_i) we reduce T by killing all the branches that are incompatible with π_i and check that the reduced tree is equivalent to C_i .

For the internship we have chosen a simple subset of the OCaml language and implemented a prototype equivalence checker along with a formal statement of correctness and proof sketches. The prototype is to be included in the OCaml compiler infrastructure and will aid the development.

2.3 From source programs to decision trees

Our source language supports integers, lists, tuples and all algebraic datatypes. Patterns support wildcards, constructors and literals, or patterns $(p_1|p_2)$ and pattern variables. We also support `when` guards. Decision trees have nodes of the form:

```

type decision_tree =
  | Unreachable
  | Failure
  | Leaf of source_expr
  | Guard of source_expr * decision_tree * decision_tree
  | Switch of accessor * (constructor * decision_tree) list * decision_tree

```

In the `Switch` node we have one subtree for every head constructor that appears in the pattern matching clauses and a fallback case for other values. The branch condition π_i expresses that the value at the switch accessor starts

with the given constructor. **Failure** nodes express match failures for values that are not matched by the source clauses. **Unreachable** is used when we statically know that no value can flow to that subtree.

We write $\llbracket t_S \rrbracket_S$ for the decision tree of the source program t_S , computed by a matrix decomposition algorithm (each column decomposition step gives a **Switch** node). It satisfies the following correctness statement:

$$\forall t_S, \forall v_S, \quad t_S(v_S) = \llbracket t_S \rrbracket_S(v_S)$$

Running any source values v_S against the source program gives the same result as running it against the decision tree.

2.4 From target programs to decision trees

The target programs include the following Lambda constructs: **let**, **if**, **switch**, **Match_failure**, **catch**, **exit**, **field** and various comparison operations, guards. The symbolic execution engine traverses the target program and builds an environment that maps variables to accessors. It branches at every control flow statement and emits a **Switch** node. The branch condition π_i is expressed as an interval set of possible values at that point. Guards result in branching. In comparison with the source decision trees, **Unreachable** nodes are never emitted.

We write $\llbracket t_T \rrbracket_T$ for the decision tree of the target program t_T , satisfying the following correctness statement:

$$\forall t_T, \forall v_T, \quad t_T(v_T) = \llbracket t_T \rrbracket_T(v_T)$$

2.5 Equivalence checking

3 Background

3.1

Objective Caml () is a dialect of the ML (Meta-Language) family of programming languages. shares many features with other dialects of ML, such as SML and Caml Light, The main features of ML languages are the use of the Hindley-Milner type system that provides many advantages with respect to static type systems of traditional imperative and object oriented language such as C, C++ and Java, such as:

- Polymorphism: in certain scenarios a function can accept more than one type for the input parameters. For example a function that com-

putes the length of a list doesn't need to inspect the type of the elements of the list and for this reason a `List.length` function can accept lists of integers, lists of strings and in general lists of any type. Such languages offer polymorphic functions through subtyping at runtime only, while other languages such as C++ offer polymorphism through compile time templates and function overloading. With the Hindley-Milner type system each well typed function can have more than one type but always has a unique best type, called the *principal type*. For example the principal type of the `List.length` function is "For any a , function from list of a to int " and a is called the *type parameter*.

- Strong typing: Languages such as C and C++ allow the programmer to operate on data without considering its type, mainly through pointers. Other languages such as C# and Go allow type erasure so at runtime the type of the data can't be queried. In the case of programming languages using an Hindley-Milner type system the programmer is not allowed to operate on data by ignoring or promoting its type.
- Type Inference: the principal type of a well formed term can be inferred without any annotation or declaration.
- Algebraic data types: types that are modelled by the use of two algebraic operations, sum and product. A sum type is a type that can hold of many different types of objects, but only one at a time. For example the sum type defined as $A + B$ can hold at any moment a value of type A or a value of type B. Sum types are also called tagged union or variants. A product type is a type constructed as a direct product of multiple types and contains at any moment one instance for every type of its operands. Product types are also called tuples or records. Algebraic data types can be recursive in their definition and can be combined.

Moreover ML languages are functional, meaning that functions are treated as first class citizens and variables are immutable, although mutable statements and imperative constructs are permitted. In addition to that features an object system, that provides inheritance, subtyping and dynamic binding, and modules, that provide a way to encapsulate definitions. Modules are checked statically and can be reified through functors.

3.2 Lambda form compilation

provides compilation in form of a bytecode executable with an optionally embeddable interpreter and a native executable that could be statically linked to provide a single file executable.

After the typechecker has proven that the program is type safe, the compiler lower the code to *Lambda*, an s-expression based language that assumes that its input has already been proved safe. On the *Lambda* representation of the source program, the compiler performs a series of optimization passes before translating the lambda form to assembly code.

1. datatypes

Most native data types in , such as integers, tuples, lists, records, can be seen as instances of the following definition

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

that is a type t with three constructors that define its complete signature. Every constructor has an arity. Nil, a constructor of arity 0, is called a constant constructor.

2. Lambda form types A lambda form target file produced by the compiler consists of a single s-expression. Every s-expression consist of (, a sequence of elements separated by a whitespace and a closing). Elements of s-expressions are:

- Atoms: sequences of ascii letters, digits or symbols
- Variables
- Strings: enclosed in double quotes and possibly escaped
- S-expressions: allowing arbitrary nesting

There are several numeric types:

- integers: that us either 31 or 63 bit two's complement arithmetic depending on system word size, and also wrapping on overflow
- 32 bit and 64 bit integers: that use 32-bit and 64-bit two's complement arithmetic with wrap on overflow
- big integers: offer integers with arbitrary precision
- floats: that use IEEE754 double-precision (64-bit) arithmetic with the addition of the literals *infinity*, *neg_infinity* and *nan*.

There are various numeric operations defined:

- Arithmetic operations: `+`, `-`, `*`, `/`, `%` (modulo), `neg` (unary negation)
- Bitwise operations: `&`, `|`, `^`, `«`, `»` (zero-shifting), `a»` (sign extending)
- Numeric comparisons: `<`, `>`, `<=`, `>=`, `==`

3. Functions

Functions are defined using the following syntax, and close over all bindings in scope: `(lambda (arg1 arg2 arg3) BODY)` and are applied using the following syntax: `(apply FUNC ARG ARG ARG)` Evaluation is eager.

4. Bindings The atom `let` introduces a sequence of bindings: `(let BINDING BINDING BINDING ... BODY)`
5. Other atoms TODO: `if`, `switch`, `stringswitch`... TODO: magari esempi

3.3 Pattern matching

Pattern matching is a widely adopted mechanism to interact with ADT. C family languages provide branching on predicates through the use of `if` statements and `switch` statements. Pattern matching on the other hand expresses predicates through syntactic templates that also allow to bind on data structures of arbitrary shapes. One common example of pattern matching is the use of regular expressions on strings. `match` provides pattern matching on ADT and primitive data types. The result of a pattern matching operation is always one of:

- this value does not match this pattern"
- this value matches this pattern, resulting the following bindings of names to values and the jump to the expression pointed at the pattern.

```
type color = | Red | Blue | Green | Black | White
```

```
match color with
| Red -> print "red"
| Blue -> print "red"
| Green -> print "red"
| _ -> print "white or black"
```

provides tokens to express data destructuring. For example we can examine the content of a list with patten matching

```
begin match list with
| [ ] -> print "empty list"
| element1 :: [ ] -> print "one element"
| (element1 :: element2) :: [ ] -> print "two elements"
| head :: tail-> print "head followed by many elements"
```

Parenthesized patterns, such as the third one in the previous example, matches the same value as the pattern without parenthesis.

The same could be done with tuples

```
begin match tuple with
| (Some _, Some _) -> print "Pair of optional types"
| (Some _, None) | (None, Some _) -> print "Pair of optional types, one of which is null"
| (None, None) -> print "Pair of optional types, both null"
```

The pattern `pattern1 | pattern2` represents the logical "or" of the two patterns `pattern1` and `pattern2`. A value matches `pattern1 | pattern2` if it matches `pattern1` or `pattern2`.

Pattern clauses can make the use of *guards* to test predicates and variables can captured (binded in scope).

```
begin match token_list with
| "switch"::var::"{":rest -> ...
| "case"::":":var::rest when is_int var -> ...
| "case"::":":var::rest when is_string var -> ...
| "}"::[ ] -> ...
| "}"::rest -> error "syntax error: " rest
```

Moreover, the pattern matching compiler emits a warning when a pattern is not exhaustive or some patterns are shadowed by precedent ones.

3.4 Symbolic execution

3.5 Translation validation

Translators, such as translators and code generators, are huge pieces of software usually consisting of multiple subsystem and constructing an actual

specification of a translator implementation for formal validation is a very long task. Moreover, different translators implement different algorithms, so the correctness proof of a translator cannot be generalized and reused to prove another translator. Translation validation is an alternative to the verification of existing translators that consists of taking the source and the target (compiled) program and proving *a posteriori* their semantic equivalence.

- Techniques for translation validation
- What does semantically equivalent mean
- What happens when there is no semantic equivalence
- Translation validation through symbolic execution

3.6 Translation validation of the Pattern Matching Compiler

1. Source program The algorithm takes as its input a source program and translates it into an algebraic data structure called *constraint_tree*.

```
type constraint_tree =
  | Unreachable
  | Failure
  | Leaf of source_expr
  | Guard of source_blackbox * constraint_tree * constraint_tree
  | Node of accessor * (constructor * constraint_tree) list * constraint_tree
```

Unreachable, Leaf of `source_expr` and Failure are the terminals of the three. We distinguish

- Unreachable: statically it is known that no value can go there
- Failure: a value matching this part results in an error
- Leaf: a value matching this part results into the evaluation of a source blackbox of code

The algorithm doesn't support type-declaration-based analysis to know the list of constructors at a given type. Let's consider some trivial examples:

```
function true -> 1
```

[] Converti a disegni

Is translated to

Node ([[true, Leaf 1]], Failure)

while

```
function
true -> 1
| false -> 2
```

will give

Node ([[true, Leaf 1]; (false, Leaf 2)])

It is possible to produce Unreachable examples by using refutation clauses (a "dot" in the right-hand-side)

```
function
true -> 1
| false -> 2
| _ -> .
```

that gets translated into Node ([[true, Leaf 1]; (false, Leaf 2)], Unreachable)

We trust this annotation, which is reasonable as the type-checker verifies that it indeed holds.

Guard nodes of the tree are emitted whenever a guard is found. Guards node contains a blackbox of code that is never evaluated and two branches, one that is taken in case the guard evaluates to true and the other one that contains the path taken when the guard evaluates to true.

[] Finisci con Node [] Spiega del fallback [] rivedi compilazione per tenere in considerazione il tuo albero invece che le lambda [] Specifica che stesso algoritmo usato per compilare a lambda, piu' optimizations

The source code of a pattern matching function in has the following form:

```

match variable with
| pattern1 -> expr1
| pattern2 when guard -> expr2
| pattern3 as var -> expr3
:
| pn -> exprn

```

and can include any expression that is legal for the compiler, such as "when" conditions and assignments. Patterns could or could not be exhaustive.

Pattern matching code could also be written using the more compact form:

```

function
| pattern1 -> expr1
| pattern2 when guard -> expr2
| pattern3 as var -> expr3
:
| pn -> exprn

```

This BNF grammar describes formally the grammar of the source program:

```

start ::= "match" id "with" patterns | "function" patterns
patterns ::= (pattern0|pattern1) pattern1+
;; pattern0 and pattern1 are needed to distinguish the first case in which
;; we can avoid writing the optional vertical line
pattern0 ::= clause
pattern1 ::= "|" clause
clause ::= lexpr "->" rexpr

lexpr ::= rule ($\varepsilon$|condition)
rexpr ::= _code ;; arbitrary code

rule ::= wildcard|variable|constructor_pattern|or_pattern ;;

;; rules
wildcard ::= "_"
variable ::= identifier

```

```

constructor_pattern ::= constructor (rule|\$\varepsilon) (assignment|\$\varepsilon)

constructor ::= int|float|char|string|bool
|unit|record|exn|objects|ref
|list|tuple|array
|variant|parameterized_variant ;; data types

or_pattern ::= wildcard|variable|constructor_pattern ("|" wildcard|variable|const

condition ::= "when" bexpr
assignment ::= "as" id
bexpr ::= _code ;; arbitrary code

```

Patterns are of the form

pattern	type of pattern
—	wildcard
x	variable
$c(p_1, p_2, \dots, p_n)$	constructor pattern
$(p_1 p_2)$	or-pattern

During compilation by the translators expressions are compiled into lambda code and are referred as lambda code actions l_i .

The entire pattern matching code is represented as a clause matrix that associates rows of patterns $(p_{i,1}, p_{i,2}, \dots, p_{i,n})$ to lambda code action l^i

$$(P \rightarrow L) = \begin{pmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n} & \rightarrow l_1 \\ p_{2,1} & p_{2,2} & \cdots & p_{2,n} & \rightarrow l_2 \\ \vdots & \vdots & \ddots & \vdots & \rightarrow \vdots \\ p_{m,1} & p_{m,2} & \cdots & p_{m,n} & \rightarrow l_m \end{pmatrix}$$

The pattern p matches a value v , written as $p \preccurlyeq v$, when one of the following rules apply

—	\preccurlyeq	v	$\forall v$
x	\preccurlyeq	v	$\forall v$
$(p_1 p_2)$	\preccurlyeq	v	iff $p_1 \preccurlyeq v$ or $p_2 \preccurlyeq v$
$c(p_1, p_2, \dots, p_a)$	\preccurlyeq	$c(v_1, v_2, \dots, v_a)$	iff $(p_1, p_2, \dots, p_a) \preccurlyeq (v_1, v_2, \dots, v_a)$
(p_1, p_2, \dots, p_a)	\preccurlyeq	(v_1, v_2, \dots, v_a)	iff $p_i \preccurlyeq v_i \forall i \in [1..a]$

When a value v matches pattern p we say that v is an *instance* of p .

Considering the pattern matrix P we say that the value vector $\vec{v} = (v_1, v_2, \dots, v_i)$ matches the line number i in P if and only if the following two conditions are satisfied:

- $p_{i,1}, p_{i,2}, \dots, p_{i,n} \preceq (v_1, v_2, \dots, v_i)$
- $\forall j < i \ p_{j,1}, p_{j,2}, \dots, p_{j,n} \not\preceq (v_1, v_2, \dots, v_i)$

We can define the following three relations with respect to patterns:

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

(a) Initial state of the compilation

Given a source of the following form:

```

match variable with
| pattern1 -> e1
| pattern2 -> e2
⋮
| pm -> em

```

the initial input of the algorithm C consists of a vector of variables $\vec{x} = (x_1, x_2, \dots, x_n)$ of size n where n is the arity of the type of x and a clause matrix $P \rightarrow L$ of width n and height m . That is:

$$C((\vec{x} = (x_1, x_2, \dots, x_n), \begin{pmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n} \rightarrow l_1 \\ p_{2,1} & p_{2,2} & \cdots & p_{2,n} \rightarrow l_2 \\ \vdots & \vdots & \ddots & \vdots \rightarrow \vdots \\ p_{m,1} & p_{m,2} & \cdots & p_{m,n} \rightarrow l_m \end{pmatrix}))$$

The base case C_0 of the algorithm is the case in which the \vec{x} is empty, that is $\vec{x} = ()$, then the result of the compilation C_0 is l_1

$$C_0((), \begin{pmatrix} \rightarrow l_1 \\ \rightarrow l_2 \\ \rightarrow \vdots \\ \rightarrow l_m \end{pmatrix})) = l_1$$

When $\vec{x} \neq ()$ then the compilation advances using one of the following four rules:

- i. Variable rule: if all patterns of the first column of P are wildcard patterns or bind the value to a variable, then

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

where

$$\begin{pmatrix} p_{1,2} & \cdots & p_{1,n} & \rightarrow & (\text{let } y_1 & x_1) & l_1 \\ p_{2,2} & \cdots & p_{2,n} & \rightarrow & (\text{let } y_2 & x_1) & l_2 \\ \vdots & \ddots & \vdots & \rightarrow & \vdots & \vdots & \vdots \\ p_{m,2} & \cdots & p_{m,n} & \rightarrow & (\text{let } y_m & x_1) & l_m \end{pmatrix}$$

That means in every lambda action l_i there is a binding of x_1 to the variable that appears on the pattern $p_{i,1}$. Bindings are omitted for wildcard patterns and the lambda action l_i remains unchanged.

- ii. Constructor rule: if all patterns in the first column of P are constructor patterns of the form $k(q_1, q_2, \dots, q_n)$ we define a new matrix, the specialized clause matrix S, by applying the following transformation on every row p :

for every $c \in \text{Set of constructors}$

for $i \leftarrow 1 \dots m$

let $k_i \leftarrow \text{constructor_of}(p_{i,1})$

if $k_i = c$ then

$P \leftarrow q_{i,1}, q_{i,2}, \dots, q_{i,n'}, p_{i,2}, p_{i,3}, \dots, p_{i,n}$

Patterns of the form $q_{i,j}$ matches on the values of the constructor and we define new fresh variables y_1, y_2, \dots, y_a so that the lambda action l_i becomes

```
(let (y1 (field 0 x1))
    (y2 (field 1 x1))
    ...
    (ya (field (a-1) x1))
    li)
```

and the result of the compilation for the set of constructors $\{c_1, c_2, \dots, c_k\}$ is:

```

switch  $x_1$  with
case  $c_1$ :  $l_1$ 
case  $c_2$ :  $l_2$ 
...
case  $c_k$ :  $l_k$ 
default: exit

```

- i. Orpat rule: there are various strategies for dealing with or-patterns. The most naive one is to split the or-patterns. For example a row p containing an or-pattern:

$$(p_{i,1} | q_{i,1} | r_{i,1}), p_{i,2}, \dots, p_{i,m} \rightarrow l_i$$

results in three rows added to the clause matrix

$$p_{i,1}, p_{i,2}, \dots, p_{i,m} \rightarrow l_i$$

$$q_{i,1}, p_{i,2}, \dots, p_{i,m} \rightarrow l_i$$

$$r_{i,1}, p_{i,2}, \dots, p_{i,m} \rightarrow l_i$$

- ii. Mixture rule: When none of the previous rules apply the clause matrix $P \rightarrow L$ is splitted into two clause matrices, the first $P_1 \rightarrow L_1$ that is the largest prefix matrix for which one of the three previous rules apply, and $P_2 \rightarrow L_2$ containing the remaining rows. The algorithm is applied to both matrices.