

Translation Verification of the pattern matching compiler

Francesco Mecca

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

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.

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

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

(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])

```

The code on 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.

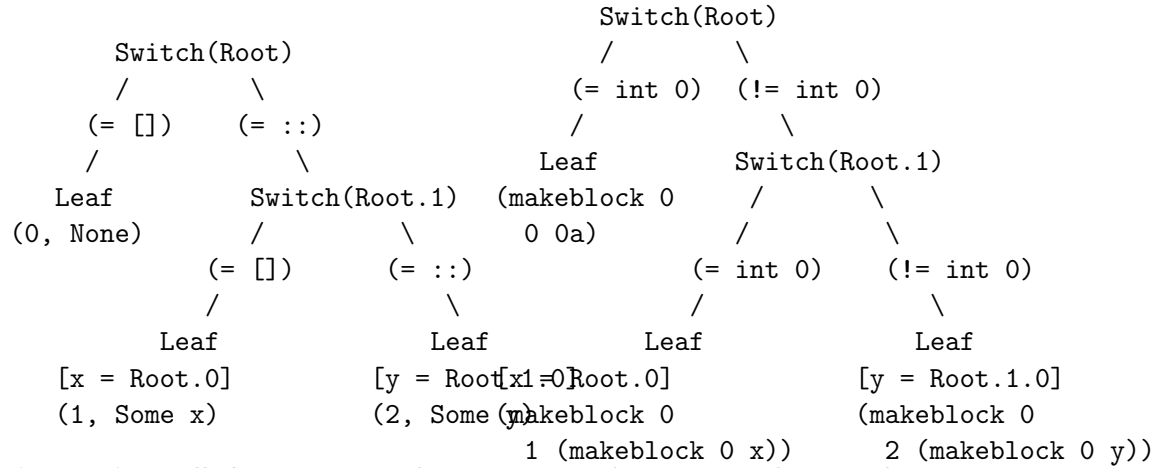
The OCaml pattern matching compiler is a critical part of the OCaml compiler in terms of correctness because bugs typically would result in wrong code production rather than triggering compilation failures. Such bugs also are hard to catch by testing because they arise in corner cases of complex patterns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

The OCaml core developers group considered evolving the pattern matching compiler, either by using a new algorithm or by incremental refactoring of its code base. 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 behavior 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 code base. The compiler is treated as a black-box and proof only depends on our equivalence algorithm.

1.1 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 are the decision trees for the source and target example program.



(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` or tuples `(x,y)` 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 .

1.2 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, which are interesting as they introduce the evaluation of expressions during matching. 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 value v_S against the source program gives the same result as running it against the decision tree.

1.3 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. In comparison with the source decision trees, **Unreachable** nodes are never emitted. 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:

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

1.4 Equivalence checking

The equivalence checking algorithm takes as input a domain of possible values S and a pair of source and target decision trees and in case the two trees are not equivalent it returns a counter example. The algorithm respects the following correctness statement:

$$\begin{aligned} \text{equiv}(S, C_S, C_T) = \text{Yes} \wedge C_T \text{ covers } S &\implies \forall v_S \approx v_T \in S, C_S(v_S) = C_T(v_T) \\ \text{equiv}(S, C_S, C_T) = \text{No}(v_S, v_T) \wedge C_T \text{ covers } S &\implies v_S \approx v_T \in S \wedge C_S(v_S) \neq C_T(v_T) \end{aligned}$$

The algorithm proceeds by case analysis. Inference rules are shown. If S is empty the results is Yes.

`equiv \emptysetset Cs Ct gs`

If the two decision trees are both terminal nodes the algorithm checks for content equality.

`equiv S Failure Failure []`

`equiv_BB BBs BBt`

`equiv S (Leaf BBs) (Leaf BBt) []`

If the source decision tree (left hand side) is a terminal while the target decision tree (right hand side) is not, the algorithm proceeds by *explosion* of the right hand side. Explosion means that every child of the right hand side is tested for equality against the left hand side.

`(equiv S Cs Ci gs)^i`
`equiv S Cs Cf gs`

`equiv S Cs (Node(a, (Domi,Ci)^i, Cf)) gs`

When the left hand side is not a terminal, the algorithm explodes the left hand side while trimming every right hand side subtree. Trimming a left hand side tree on an interval set dom_S computed from the right hand side tree constructor means mapping every branch condition dom_T (interval set of possible values) on the left to the intersection of dom_T and dom_S when the accessors on both side are equal, and removing the branches that result in an empty intersection. If the accessors are different, dom_T is left unchanged.

`equiv S Ci (trim Ct a=Ki) gs`
`equiv S Cf (trim Ct (a\notin(Ki)^i) gs`

`equiv S (Node(a, (Ki,Ci)^i, Cf) Ct gs`

The equivalence checking algorithm deals with guards by storing a queue. A guard blackbox is pushed to the queue whenever the algorithm encounters

a Guard node on the right, while it pops a blackbox from the queue whenever a Guard node appears on the left hand side. The algorithm stops with failure if the popped blackbox and the and blackbox on the left hand Guard node are different, otherwise it continues by exploding to two subtrees, one in which the guard condition evaluates to true, the other when it evaluates to false. Termination of the algorithm is successful only when the guards queue is empty.

```

equiv S Ctrue Ct (gs++[condition])
equiv S Cfalse Ct (gs++[condition])
-----
equiv S (Guard condition Ctrue Cfalse) Ct gs

equiv S Cs Ctrue gs
equiv S Cs Cfalse gs
-----
equiv S Cs (Guard condition Ctrue Cfalse) ([condition]++gs)

```

2 Background

2.1 OCaml

Objective Caml (OCaml) is a dialect of the ML (Meta-Language) family of programming that 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 computes 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 modeled 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.

2.2 Compiling OCaml code

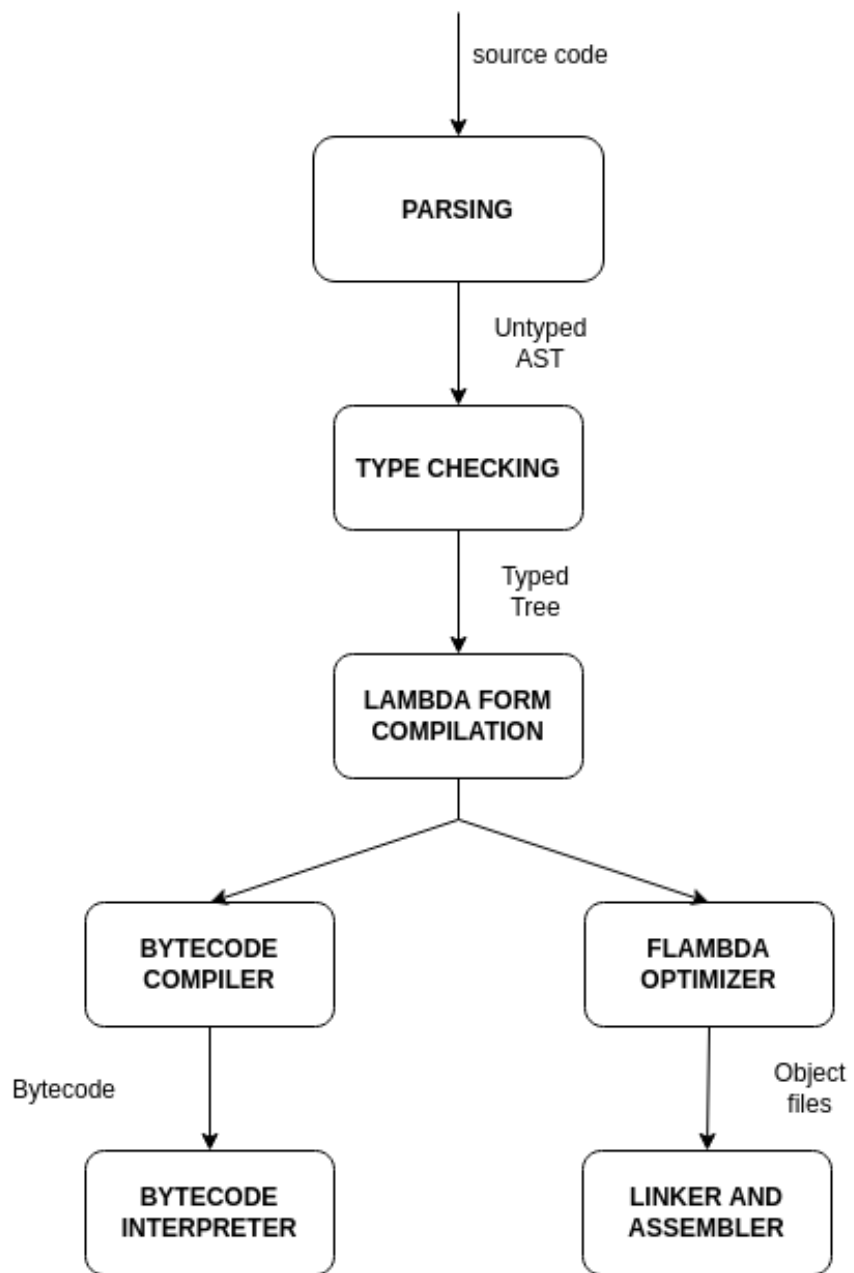
The OCaml compiler provides compilation of source files in form of a byte-code executable with an optionally embeddable interpreter or as a native executable that could be statically linked to provide a single file executable. Every source file is treated as a separate *compilation unit* that is advanced through different states. The first stage of compilation is the parsing of the input code that is transformed into an untyped syntax tree. Code with syntax errors is rejected at this stage. After that the AST is processed by the type checker that performs three steps at once:

- type inference, using the classical *Algorithm W*
- perform subtyping and gathers type information from the module system

- ensures that the code obeys the rule of the OCaml type system

At this stage, incorrectly typed code is rejected. In case of success, the untyped AST is transformed into a *Typed Tree*. After the typechecker has proven that the program is type safe, the compiler lowers the code to *Lambda*, an s-expression based language that assumes that its input has already been proved safe. After the Lambda pass, the Lambda code is either translated into bytecode or goes through a series of optimization steps performed by the *Flambda* optimizer before being translated into assembly.

This is an overview of the different compiler steps.



2.3 Memory representation of OCaml values

An usual OCaml source program contains few to none explicit type signatures. This is possible because of type inference that allows to annotate the AST with type informations. However, since the OCaml typechecker guarantees that a program is well typed before being transformed into Lambda code, values at runtime contains only a minimal subset of type informations needed to distinguish polymorphic values. For runtime values, OCaml uses a uniform memory representation in which every variable is stored as a value in a contiguous block of memory. Every value is a single word that is either a concrete integer or a pointer to another block of memory, that is called *cell* or *box*. We can abstract the type of OCaml runtime values as the following:

```
type t = Constant | Cell of int * t
```

where a one bit tag is used to distinguish between Constant or Cell. In particular this bit of metadata is stored as the lowest bit of a memory block.

Given that all the OCaml target architectures guarantee that all pointers are divisible by four and that means that two lowest bits are always 00 storing this bit of metadata at the lowest bit allows an optimization. Constant values in OCaml, such as integers, empty lists, Unit values and constructors of arity zero (*constant* constructors) are unboxed at runtime while pointers are recognized by the lowest bit set to 0.

2.4 Lambda form compilation

A Lambda code target file is produced by the compiler and 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

The Lambda form is a key stage where the compiler discards type informations and maps the original source code to the runtime memory model described. In this stage of the compiler pipeline pattern match statements are analyzed and compiled into an automata.

```
type t = | Foo | Bar | Baz | Fred
```

```
let test = function  
  | Foo -> "foo"  
  | Bar -> "bar"  
  | Baz -> "baz"  
  | Fred -> "fred"
```

The Lambda output for this code can be obtained by running the compiler with the *-dlambda* flag:

```
(setglobal Prova!  
  (let  
    (test/85 =  
      (function param/86  
        (switch* param/86  
          case int 0: "foo"  
          case int 1: "bar"  
          case int 2: "baz"  
          case int 3: "fred"))))  
    (makeblock 0 test/85)))
```

As outlined by the example, the *makeblock* directive is responsible for allocating low level OCaml values and every constant constructor of the algebraic type *t* is stored in memory as an integer. The *setglobal* directives declares a new binding in the global scope: Every concept of modules is lost at this stage of compilation. The pattern matching compiler uses a jump table to map every pattern matching clauses to its target expression. Values are addressed by a unique name.

```
type t = | English of p | French of q  
type p = | Foo | Bar  
type q = | Tata| Titi  
type t = | English of p | French of q
```

```
let test = function  
  | English Foo -> "foo"  
  | English Bar -> "bar"  
  | French Tata -> "baz"  
  | French Titi -> "fred"
```

In the case of types with a smaller number of variants, the pattern matching compiler may avoid the overhead of computing a jump table. This example

also highlights the fact that non constant constructor are mapped to cons cell that are accessed using the *tag* directive.

```
(setglobal Prova!  
  (let  
    (test/89 =  
      (function param/90  
(switch* param/90  
  case tag 0: (if (≠ (field 0 param/90) 0) "bar" "foo")  
  case tag 1: (if (≠ (field 0 param/90) 0) "fred" "baz"))))  
    (makeblock 0 test/89)))
```

In the Lambda language 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*.

The are various numeric operations defined:

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

1. 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.

2. Other atoms

The atom *let* introduces a sequence of bindings at a smaller scope than the global one: (let BINDING BINDING BINDING ... BODY)

The Lambda form supports many other directives such as *strinswitch* that is constructs aspecialized jump tables for string, integer range

comparisons and so on. These constructs are explicitly undocumented because the Lambda code intermediate language can change across compiler releases.

2.5 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. 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 pattern 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 be captured (bound 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.

2.6 Symbolic execution

Symbolic execution is a widely used technique in the field of computer security. It allows to analyze different execution paths of a program simultaneously while tracking which inputs trigger the execution of different parts of the program. Inputs are modelled symbolically rather than taking "concrete" values. A symbolic execution engine keeps track of expressions and variables in terms of these symbolic symbols and attaches logical constraints to every branch that is being followed. Symbolic execution engines are used to track bugs by modelling the domain of all possible inputs of a program, detecting infeasible paths, dead code and proving that two code segments are equivalent.

Let's take as example this signedness bug that was found in the FreeBSD kernel and allowed, when calling the `getpeername` function, to read portions of kernel memory.

```
int compat;
```

```

{
    struct file *fp;
    register struct socket *so;
    struct sockaddr *sa;
    int len, error;

    ...

    len = MIN(len, sa->sa_len);    /* [1] */
    error = copyout(sa, (caddr_t)uap->asa, (u_int)len);
    if (error)
goto bad;

    ...

bad:
    if (sa)
FREE(sa, M_SONAME);
    fdrop(fp, p);
    return (error);
}

```

The tree of the execution when the function is evaluated considering *int len* our symbolic variable α , *sa->sa_len* as symbolic variable β and π as the set of constraints on a symbolic variable:

```

[1]          compat (...)          {  $\pi_{\alpha}$ :  $-\infty < \alpha < \infty$  }
|
[2]          min ( $\sigma_1$ ,  $\sigma_2$ )      {  $\pi_{\sigma}$ :  $-\infty < \sigma < \infty$  }
|
[3]          cast(u_int) (...)      {  $\pi_{\sigma}$ :  $0 \leq \sigma < \infty$  }
|
... // rest of the execution

```

We can see that at step 3 the set of possible values of the scrutinee α is bigger than the set of possible values of the input σ to the *cast* directive, that is: $\pi_{\alpha} \not\subseteq \pi_{\sigma}$. For this reason the *cast* may fail when α is *len* negative number, outside the domain π_{σ} . In C this would trigger undefined behaviour (signed overflow) that made the exploitation possible.

Every step of evaluation can be modelled as the following transition:

$$(\pi_{\sigma}, \pi_i) \rightarrow (\pi'_{\sigma}, \pi'_i)$$

if we express the π constraints as logical formulas we can model the execution of the program in terms of Hoare Logic. State of the computation is a Hoare triple $\{P\}C\{Q\}$ where P and Q are respectively the *precondition* and the *postcondition* that constitute the assertions of the program. C is the directive being executed. The language of the assertions P is:

$$P ::= \text{true} \mid \text{false} \mid a < b \mid P_1 \wedge P_2 \mid P_1 \vee P_2 \mid \sim P$$

where a and b are numbers. In the Hoare rules assertions could also take the form of

$$P ::= \forall i. P \mid \exists i. P \mid P_1 \rightarrow P_2$$

where i is a logical variable, but assertions of these kinds increases the complexity of the symbolic engine. Execution follows the rules of Hoare logic:

- Empty statement :

$$\text{-----}$$

$$\{P\}/\text{skip}/\{P\}$$

- Assignment statement : The truthness of $P[a/x]$ is equivalent to the truth of $\{P\}$ after the assignment.

$$\text{-----}$$

$$\{P[a/x]\}x:=a\{P\}$$

- Composition : c_1 and c_2 are directives that are executed in order; $\{Q\}$ is the *midcondition*.

$$\{P\}c_1\{R\}, \{R\}c_2\{Q\}$$

$$\text{-----}$$

$$\{P\}c_1;c_2\{Q\}$$

- Conditional :

$$\{P \wedge b\}c_1\{Q\}, \{P \wedge \sim b\}c_2\{Q\}$$

$$\text{-----}$$

$$\{P\}\text{if } b \text{ then } c_1 \text{ else } c_2\{Q\}$$

- Loop : $\{P\}$ is the loop invariant. After the loop is finished P holds and $\sim b$ caused the loop to end.

$\{P \wedge b\}c\{P\}$

 $\{P\}\text{while } b \text{ do } c\{P \wedge \sim b\}$

Even if the semantics of symbolic execution engines are well defined, the user may run into different complications when applying such analysis to non trivial codebases. For example, depending on the domain, loop termination is not guaranteed. Even when termination is guaranteed, looping causes exponential branching that may lead to path explosion or state explosion. Reasoning about all possible executions of a program is not always feasible and in case of explosion usually symbolic execution engines implement heuristics to reduce the size of the search space.

2.7 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 Translation validation of the Pattern Matching Compiler

3.1 Source program

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

```

type decision_tree =
  | Unreachable
  | Failure
  | Leaf of source_expr
  | Guard of source_blackbox * decision_tree * decision_tree
  | Node of accessor * (constructor * decision_tree) list * decision_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 black box 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 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 OCaml compiler, such as *when* guards 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 ::= rule ("|" wildcard|variable|constructor_pattern)+

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 $p \preceq v$, when one of the following rules apply

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

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:

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

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

3.1.1 Parsing of the source program

The source program of the following general form is parsed using a parser generated by Menhir, a LR(1) parser generator for the OCaml programming language. Menhir compiles LR(1) a grammar specification, in this case the OCaml pattern matching grammar, down to OCaml code.

```

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

```

The result of parsing, when successful, results in a list of clauses and a list of type declarations. Every clause consists of three objects: a left-hand-side that is the kind of pattern expressed, an option guard and a right-hand-side expression. Patterns are encoded in the following way:

pattern	type
-	Wildcard
p_1 as x	Assignment
$c(p_1, p_2, \dots, p_n)$	Constructor
$(p_1 \mid p_2)$	Orpat

Guards and right-hand-sides are treated as a blackbox of OCaml code. A sound approach for treating these blackbox would be to inspect the OCaml compiler during translation to Lambda code and extract the blackboxes compiled in their Lambda representation. This would allow to test for equality with the respective blackbox at the target level. Given that this level of introspection is currently not possible, we decided to restrict the structure of blackboxes to the following (valid) OCaml code:

```
external guard : 'a -> 'b = "guard"
external observe : 'a -> 'b = "observe"
```

We assume these two external functions *guard* and *observe* with a valid type that lets the user pass any number of arguments to them. All the guards are of the form `guard <arg> <arg> <arg>`, where the `<arg>` are expressed using the OCaml pattern matching language. Similarly, all the right-hand-side expressions are of the form `observe <arg> <arg> ...` with the same constraints on arguments.

```
type t = Z | S of t
```

```
let _ = function
  | Z -> observe 0
  | S Z -> observe 1
  | S x when guard x -> observe 2
  | S (S x) as y when guard x y -> observe 3
  | S _ -> observe 4
```

Once parsed, the type declarations and the list of clauses are encoded in the form of a matrix that is later evaluated using a matrix decomposition algorithm.

3.1.2 Matrix decomposition of pattern clauses

The initial input of the decomposition 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

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

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

1. 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

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

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.

2. Constructor rule: if all patterns in the first column of P are constructors 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 x1 with
case c1: l1
case c2: l2
...
case ck: lk
default: exit
```

1. 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

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

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

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

2. Mixture rule: When none of the previous rules apply the clause matrix $P \rightarrow L$ is split 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.

4 Correctness of the algorithm

Running a program t_s or its translation $\llbracket t_s rbracket$ against an input v_s produces as a result a result r in the following way:

$$\begin{aligned} & (\llbracket t_s rbracket(v_s) = C_s(v_s)) \rightarrow r \\ & t_s(v_s) \rightarrow r \end{aligned}$$

Likewise

$$\begin{aligned} & (\llbracket t_t rbracket(v_t) = C_t(v_t)) \rightarrow r \\ & t_t(v_t) \rightarrow r \end{aligned}$$

where result $r ::= \text{guard list} * (\text{Match blackbox} \mid \text{NoMatch} \mid \text{Absurd})$ and $\text{guard} ::= \text{blackbox}$.

Having defined equivalence between two inputs of which one is expressed in the source language and the other in the target language $v_s \simeq v_t$ (TODO define, this talks about the representation of source values in the target)

we can define the equivalence between a couple of programs or a couple of decision trees

$$\begin{aligned} t_s \simeq t_t & := \forall v_s \simeq v_t, t_s(v_s) = t_t(v_t) \\ C_s \simeq C_t & := \forall v_s \simeq v_t, C_s(v_s) = C_t(v_t) \end{aligned}$$

The proposed equivalence algorithm that works on a couple of decision trees is returns either *Yes* or *No*(v_s, v_t) where v_s and v_t are a couple of possible counter examples for which the constraint trees produce a different result.

4.1 Statements

Theorem. We say that a translation of a source program to a decision tree is correct when for every possible input, the source program and its respective decision tree produces the same result

$$\forall v_s, t_s(v_s) = \llbracket t_s rbracket(v_s)$$

Likewise, for the target language:

$$\forall v_t, t_t(v_t) = \llbracket t_t rbracket(v_t)$$

Definition: in the presence of guards we can say that two results are equivalent modulo the guards queue, written $r_1 \simeq_{gs} r_2$, when:

$$(gs_1, r_1) \simeq_{gs} (gs_2, r_2) \Leftrightarrow (gs_1, r_1) = (gs_2 ++ gs, r_2)$$

Definition: we say that C_t covers the input space S , written $\text{/covers}(C_t, S)$ when every value $v_s \in S$ is a valid input to the decision tree C_t . (TODO: rephrase)

Theorem: Given an input space S and a couple of decision trees, where the target decision tree C_t covers the input space S , we say that the two decision trees are equivalent when:

$$\text{equiv}(S, C_s, C_t, gs) = \text{Yes} \wedge \text{covers}(C_t, S) \rightarrow \forall v_s \simeq v_t \in S, C_s(v_s) \simeq_{gs} C_t(v_t)$$

Similarly we say that a couple of decision trees in the presence of an input space S are *not* equivalent when:

$$\text{equiv}(S, C_s, C_t, gs) = \text{No}(v_s, v_t) \wedge \text{covers}(C_t, S) \rightarrow v_s \simeq v_t \in S \wedge C_s(v_s) \neq_{gs} C_t(v_t)$$

Corollary: For a full input space S , that is the universe of the target program we say:

$$\text{equiv}(S, \llbracket t_s rbracket_s, \llbracket t_t rbracket_t, \emptyset) = \text{Yes} \Leftrightarrow t_s \simeq t_t$$

1. Proof of the correctness of the translation from source programs to source decision trees

We define a source term t_s as a collection of patterns pointing to black-boxes

$$t_s ::= (p \rightarrow \text{bb})^{i \in I}$$

A pattern is defined as either a constructor pattern, an or pattern or a constant pattern

$$p ::= K(p_i)^i, i \in I \quad (p \quad q) \quad n \in \mathbb{N}$$

A decision tree is defined as either a Leaf, a Failure terminal or an intermediate node with different children sharing the same accessor a and an optional fallback. Failure is emitted only when the patterns don't cover the whole set of possible input values S . The fallback is not needed when the user doesn't use a wildcard pattern.

$$\begin{aligned}
C_s &::= \text{Leaf } bb \quad \text{Node}(a, (K_i \rightarrow C_i)^{i \in S}, C?) \\
a &::= \text{Here} \quad n.a \\
v_s &::= K(v_i)^{i \in I} \quad n \in \mathbb{N}
\end{aligned}$$

We define the decomposition matrix m_s as

$$\text{SMatrix } m_s := (a_j)^{j \in J}, ((p_{ij})^{j \in J} \rightarrow bb_i)^{i \in I}$$

We define the decision tree of source programs $[[t_s]]$ in terms of the decision tree of pattern matrices $[[m_s]]$ by the following: $[[((p_i \rightarrow bb_i)^{i \in I} rbracket := [[(\text{Root}), (p_i \rightarrow bb_i)^{i \in I} rbracket \# i \in I$

decision tree computed from pattern matrices respect the following invariant:

$$\forall v (v_i)^{i \in I} = v(a_i)^{i \in I} \rightarrow [[mrbracket(v) = m(v_i)^{i \in I} \text{ for } m = ((a_i)^{i \in I}, (r_i)^{i \in I})$$

where

$$\begin{aligned}
v(\text{Here}) &= v \\
K(v_i)^i(k.a) &= v_k(a) \text{ if } k \in [0;n[
\end{aligned}$$

We proceed to show the correctness of the invariant by a case analysys.

Base cases:

- (a) $[[\emptyset, (\emptyset \rightarrow bb_i)^i]] := \text{Leaf } bb_i$ where $i := \min(I)$, that is a decision tree $[[m_s]]$ defined by an empty accessor and empty patterns pointing to blackboxes bb_i . This respects the invariant because a decomposition matrix in the case of empty rows returns the first expression and we known that $(\text{Leaf } bb)(v) := \text{Match } bb$
- (b) $[[(a_j)^j, \emptyset]]$:= Failure

Regarding non base cases: We define

$$\begin{aligned}
\text{let } \text{Idx}(k) &= [0; \text{arity}(k)[\\
\text{let } \text{First}(\emptyset) &= \perp \\
\text{let } \text{First}((a_j)^j) &= a_{\min(j)} \# \# \text{ where } j \in J \neq \emptyset
\end{aligned}$$

4.2 Proof of equivalence checking

The equivalence checking algorithm takes as parameters an input space S , a source decision tree C_s and a target decision tree C_t :

$$\text{equiv}(S, C_s, C_t) \rightarrow \text{Yes} \quad \text{No}(v_s, v_t)$$

When the algorithm returns Yes and the input space is covered by C_s we can say that the couple of decision trees are the same for every couple of source value v_s and target value v_t that are equivalent.

$$\text{equiv}(S, C_s, C_t) = \text{Yes and cover}(C_t, S) \rightarrow \forall v_s \simeq v_t \in S \wedge C_s(v_s) = C_t(v_t)$$

In the case where the algorithm returns No we have at least a couple of counter example values v_s and v_t for which the two decision trees outputs a different result.

$$\text{equiv}(S, C_s, C_t) = \text{No}(v_s, v_t) \text{ and cover}(C_t, S) \rightarrow \forall v_s \simeq v_t \in S \wedge C_s(v_s) \neq C_t(v_t)$$

We define the following

$$\begin{aligned} \text{Forall}(\text{Yes}) &= \text{Yes} \\ \text{Forall}(\text{Yes}::l) &= \text{Forall}(l) \\ \text{Forall}(\text{No}(v_s, v_t)::_) &= \text{No}(v_s, v_t) \end{aligned}$$

There exists and are injective:

$$\begin{aligned} \text{int}(k) &\in \mathbb{N} \quad (\text{ar}(k) = 0) \\ \text{tag}(k) &\in \mathbb{N} \quad (\text{ar}(k) > 0) \\ \pi(k) &= \{n \\ \text{int}(k) = n\} \times \{n \\ \text{tag}\{k\} = n\} \end{aligned}$$

where k is a constructor.

We proceed by case analysis:

1. $\text{equiv}(\emptyset, C_s, C_t) := \text{Yes}$

In the other subcases S is always non-empty.

1. $\text{equiv}(S, \text{Failure}, \text{Failure}) := \text{Yes}$ the statement holds because of equality between Failure nodes in the case of every possible value v .

2. The result of the subcase where we have a source decision tree C_s that is either a Leaf terminal or a Failure terminal and a target decision tree defined by an accessor a and a positive number of couples constraint π_i and children nodes C_{ti} . The output the output of the algorithm is:

$$\text{equiv}(S, (\text{Leaf } \text{bb}_s \quad \text{Failure}) \text{ as } C_s, \text{Node}(a, (\pi_i \rightarrow C_{ti})^i)) := \text{Forall}(\text{equiv}(S_{a \rightarrow \pi(k_i)}), C_s, C_{ti})$$

The statement holds because defined let $S_i := S(a \rightarrow \pi_i)$ either the algorithm is true for every sub-input space S_i and subtree C_{ti}

$$\text{equiv}(S_i, C_s, C_{ti}) = \text{Yes } \forall i$$

or we have a counter example v_s, v_t for which

$$v_s \simeq v_t \in S_k \wedge c_s(v_s) \neq C_{tk}(v_t)$$

then because

$$v_t \in (a \rightarrow \pi_k) \Rightarrow C_t(v_t) = C_{tk}(v_t)$$

then

$$v_s \simeq v_t \in S \wedge C_s(v_s) \neq C_t(v_t)$$

and the result of the algorithm is

$$\text{equiv}(S_i, C_s, C_{ti}) = \text{No}(v_s, v_t) \text{ for some minimal } k \in I$$