

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

Given as input to the pattern matching compiler, this snippet of code gets translated into the Lambda intermediate representation of the OCaml com-

piler. The Lambda representation of a program is shown by calling the `ocamlc` compiler with `-drawlambda` flag. In this example we renamed the variables assigned in order to ease the understanding of the tests that are performed when the code is translated into the Lambda form. code phase.

```
(function scrutinee
  (if scrutinee ;;; true when scrutinee (list) not empty
    (let (tail =a (field 1 scrutinee/81)) ;;; assignment
      (if tail
        (let
          y =a (field 0 tail))
          ;;; y is the first element of the tail
          (makeblock 0 2 (makeblock 0 y)))
          ;;; allocate memory for tuple (2, Some y)
        (let (x =a (field 0 scrutinee))
          ;;; x is the head of the scrutinee
          (makeblock 0 1 (makeblock 0 x))))
          ;;; allocate memory for tuple (1, Some x)
        [0: 0 0a])) ;;; low level representatio of (0, None)
```

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.

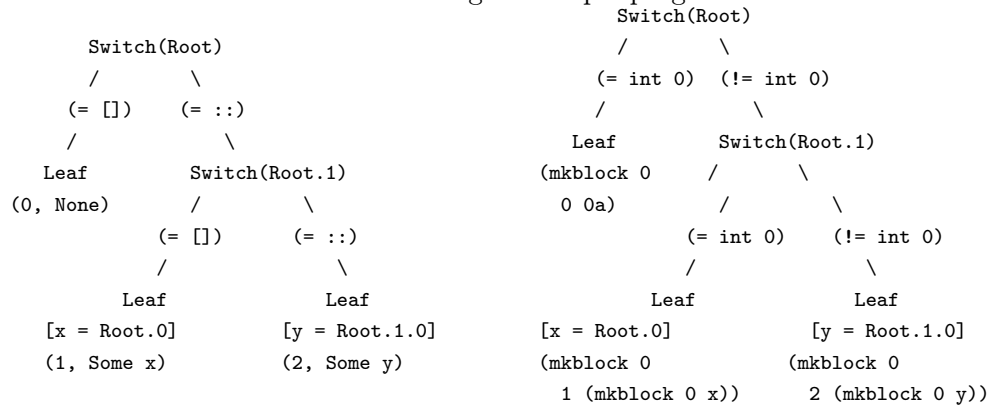
The OCaml core developers group considered evolving the pattern matching compiler, either by using a new algorithm or by incremental refactoring of its code base. For this reason we want to verify that new implementations of the compiler avoid the introduction of new bugs and that such modifications 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

Our algorithm translates both source and target programs into a common representation, decision trees. Decision trees where 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  $\pi_i$  and a subtree  $C_i$ . For every child  $(\pi_i, C_i)$  we reduce  $T$  by killing all the branches that are incompatible with  $\pi_i$  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 such as  $(p_1|p_2)$  and pattern variables. In particular Or-patterns provide a more compact way to group patterns that point to the same expression.

```

match w with
| p1 -> expr           match w with
| p2 -> expr           | p1|p2|p3 -> expr
| p3 -> expr

```

We also support **when** guards, which are interesting as they introduce the evaluation of expressions during matching. This is the type definition of decision tree as they are used in the prototype implementation:

```

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  $\pi_i$  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  $\llbracket t_s \rrbracket_s$  to denote the translation of the source program (the set of pattern matching clauses) into a decision tree, computed by a matrix decomposition algorithm (each column decomposition step gives a **Switch** node). It satisfies the following correctness statement:

$$\forall t_s, \forall v_s, \quad t_s(v_s) = \llbracket t_s \rrbracket_s(v_s)$$

The correctness statement intuitively states that for every source program, for every source value that is well-formed input to a source program, running the program  $t_S$  against the input value  $v_S$  is the same as running the compiled source program  $\llbracket t_S \rrbracket$  (that is a decision tree) against the same input value  $v_S$ .

### 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  $\pi_i$  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  $\llbracket t_T \rrbracket_T$  to denote the translation of a target program  $t_T$  into a decision tree of the target program  $t_T$ , satisfying the following correctness statement that is symmetric to the correctness statement for the translation of source programs:

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

### 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. Our 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.

$$\frac{}{\text{equiv}(\emptyset, C_S, C_T)G}$$

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

$$\frac{}{\text{equiv}(S, \text{Failure}, \text{Failure})[]}$$

$$\frac{t_S \approx_{\text{term}} t_T}{\text{equiv}(S, \text{Leaf}(t_S), \text{Leaf}(t_T))[]}$$

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.

$$\frac{C_S \in \text{Leaf}(t), \text{Failure} \quad \forall i, \text{equiv}((S \wedge a \in D_i), C_S, C_i)G \quad \text{equiv}((S \wedge a \notin (D_i)^i), C_S, C_{\text{fb}})G}{\text{equiv}(S, C_S, \text{Switch}(a, (D_i)^i C_i, C_{\text{fb}}))G}$$

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.

$$\frac{\forall i, \text{equiv}((S \wedge a = K_i), C_i, \text{trim}(C_T, a = K_i))G \quad \text{equiv}((S \wedge a \notin (K_i)^i), C_{\text{fb}}, \text{trim}(C_T, a \notin (K_i)^i))G}{\text{equiv}(S, \text{Switch}(a, (K_i, C_i)^i, C_{\text{fb}}), C_T)G}$$

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 in 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.

$$\frac{\text{equiv}(S, C_0, C_T)G, (t_S = 0) \quad \text{equiv}(S, C_1, C_T)G, (t_S = 1)}{\text{equiv}(S, \text{Guard}(t_S, C_0, C_1), C_T)G}$$

$$\frac{t_S \approx_{\text{term}} t_T \quad \text{equiv}(S, C_S, C_b)G}{\text{equiv}(S, C_S, \text{Guard}(t_T, C_0, C_1))(t_S = b), G}$$

## 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.



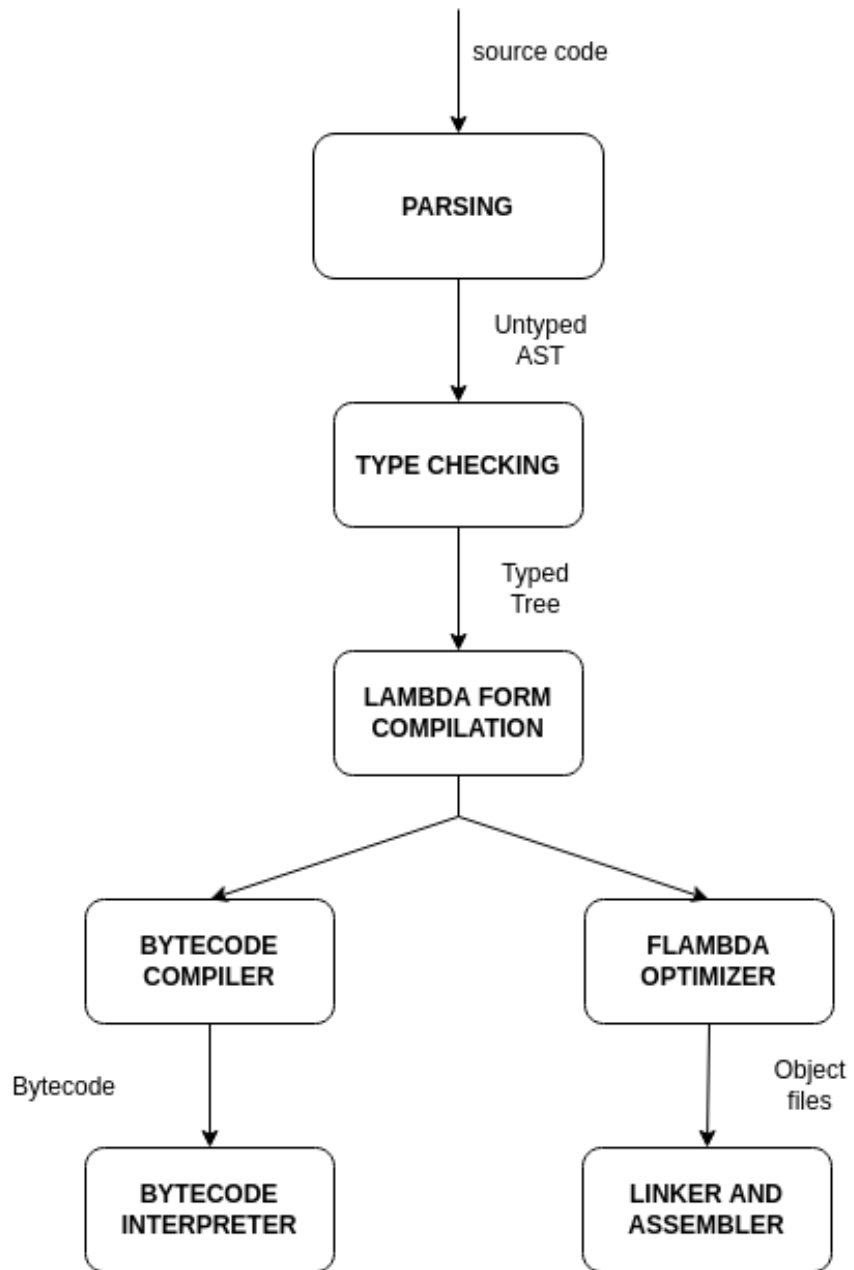
## 2.2 Compiling OCaml code

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 information 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*.

There 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 constructs specialized 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 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.

## 2.6 Symbolic execution

Symbolic execution is a widely used techniques 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;

    ...

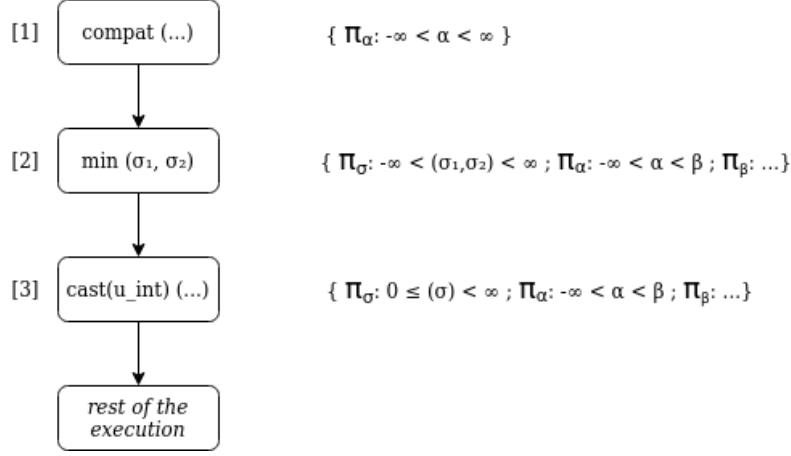
    len = MIN(len, sa->sa_len);    /* [1] */
    error = copyout(sa, (caddr_t)uap->asa, (u_int)len);
    if (error)
goto bad;

    ...

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  $\alpha$ , *sa->sa\_len* as symbolic variable  $\beta$  and  $\pi$  as the set of constraints on a symbolic variable:



We can see that at step 3 the set of possible values of the scrutinee  $\alpha$  is bigger than the set of possible values of the input  $\sigma$  to the *cast* directive, that is:  $\pi_\alpha \not\subseteq \pi_\sigma$ . For this reason the *cast* may fail when  $\alpha$  is *len* negative number, outside the domain  $\pi_\sigma$ . 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:

$$(\pi_\sigma, (\pi_i)^i) \rightarrow (\pi'_\sigma, (\pi'_i)^i)$$

if we express the  $\pi$  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:

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

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 :

$$\overline{\{P\}skip\{P\}}$$

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

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

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

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

- Conditional :

$$\frac{\{P \wedge b\}c_1\{Q\}, \{P \wedge \neg b\}c_2\{Q\}}{\{P\}if\ b\ then\ c_1\ else\ c_2\{Q\}}$$

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

$$\frac{\{P \wedge b\}c\{P\}}{\{P\}while\ b\ do\ c\{P \wedge \neg b\}}$$

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

Our algorithm takes as its input a source program and translates it into an algebraic data structure which type we call *decision\_tree*.

```
type decision_tree =  
  | Unreachable  
  | Failure  
  | Leaf of source_expr  
  | Guard of source_blackbox * decision_tree * decision_tree  
  | Switch 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

Our 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
```

is translated to

```
Switch ([[true, Leaf 1]], Failure)
```

while

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

will be translated to

```
Switch ([[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

```
Switch ([[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.

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
```

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 ;;
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|parameter
or _pattern ::= rule ("|" wildcard|variable|constructor _pattern)+
condition ::= "when" bexpr
assignment ::= "as" id
bexpr ::= _code ;; arbitrary code

```

The source program is parsed using the `ocaml-compiler-libs` library. 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
<code>_</code>	Wildcard
<code>p<sub>1</sub> as x</code>	Assignment
<code>c(p<sub>1</sub>,p<sub>2</sub>,... ,p<sub>n</sub>)</code>	Constructor
<code>(p<sub>1</sub>   p<sub>2</sub>)</code>	Orpat

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.

Patterns are of the form

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

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

–	$\preccurlyeq$	$v$	$\forall v$
x	$\preccurlyeq$	$v$	$\forall v$
$(p_1 \mid p_2)$	$\preccurlyeq$	$v$	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  $p$ .

During compilation by the translators, expressions at the right-hand-side are compiled into Lambda code and are referred as lambda code actions  $l_i$ .

We define the *pattern matrix*  $P$  as the matrix  $|m \times n|$  where  $m$  bigger or equal than the number of clauses in the source program and  $n$  is equal to the arity of the constructor with the gratest arity.

$$P = \begin{pmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n} \\ p_{2,1} & p_{2,2} & \cdots & p_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m,1} & p_{m,2} & \cdots & p_{m,n} \end{pmatrix}$$

every row of  $P$  is called a pattern vector  $\vec{p}_i = (p_1, p_2, \dots, p_n)$ ; In every instance of  $P$  pattern vectors appear normalized on the length of the longest pattern vector. Considering the pattern matrix  $P$  we say that the value vector  $\vec{v} = (v_1, v_2, \dots, v_i)$  matches the pattern vector  $p_i$  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

With the support of two auxiliary functions

- tail of an ordered family

$$\text{tail}((x_i)^{i \in I}) := (x_i)^{i \neq \min(I)}$$

- first non- $\perp$  element of an ordered family

$$\begin{aligned} \text{First}((x_i)^i) &:= \perp && \text{if } \forall i, x_i = \perp \\ \text{First}((x_i)^i) &:= x_{\min\{i \mid x_i \neq \perp\}} && \text{if } \exists i, x_i \neq \perp \end{aligned}$$

we now define what it means to run a pattern row against a value vector of the same length, that is  $(p_i)^i(v_i)^i$

$p_i$	$v_i$	$\text{result}_{\text{pat}}$
$\emptyset$	$(\emptyset)$	$\llbracket$
$(\_, \text{tail}(p_i)^i)$	$(v_i)$	$\text{tail}(p_i)^i(\text{tail}(v_i)^i)$
$(x, \text{tail}(p_i)^i)$	$(v_i)$	$\sigma[x \mapsto v_0]$ if $\text{tail}(p_i)^i(\text{tail}(v_i)^i) = \sigma$
$(K(q_j)^j, \text{tail}(p_i)^i)$	$(K(v'_j)^j, \text{tail}(v_j)^j)$	$((q_j)^j + \text{tail}(p_i)^i)((v'_j)^j + \text{tail}(v_i)^i)$
$(K(q_j)^j, \text{tail}(p_i)^i)$	$(K'(v'_l)^l, \text{tail}(v_j)^j)$	$\perp$ if $K \neq K'$
$(q_1 q_2, \text{tail}(p_i)^i)$	$(v_i)^i$	$\text{First}((q_1, \text{tail}(p_i)^i)(v_i)^i, (q_2, \text{tail}(p_i)^i)(v_i)^i)$

A source program  $t_S$  is a collection of pattern clauses pointing to *bb* terms. Running a program  $t_S$  against an input value  $v_S$ , written  $t_S(v_S)$  produces a result  $r$  belonging to the following grammar:

$$\begin{aligned} t_S &::= (p \rightarrow bb)^{i \in I} \\ t_S(v_S) &\rightarrow r \\ r &::= \text{guard list} * (\text{Match}(\sigma, bb) \mid \text{NoMatch} \mid \text{Absurd}) \end{aligned}$$

We can define what it means to run an input value  $v_S$  against a source program  $t_S$ :

$$\begin{aligned} \text{First}((x_i)^i) &:= x_{\min\{i \mid x_i \neq \perp\}} \quad \text{if } \exists i, x_i \neq \perp \\ t_S(v_S) &:= \text{NoMatch} \quad \text{if } \forall i, p_i(v_S) = \perp \\ t_S(v_S) &= \{ \text{Absurd if } bb_{i_o} = . \text{ (refutation clause)} \\ &\quad \text{Match}(\sigma, bb_{i_o}) \text{ otherwise} \\ &\quad \text{where } i_o = \min\{i \mid p_i(v_S) \neq \perp\} \end{aligned}$$

*guard* and *bb* expressions are treated as blackboxes of OCaml code. A sound approach for treating these blackboxes 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.

```

(* declaration of an algebraic and recursive datatype t *)
type t = K1 | K2 of t

let _ = function
  | K1 -> observe 0
  | K2 K1 -> observe 1
  | K2 x when guard x -> observe 2 (* guard inspects the x variable *)
  | K2 (K2 x) as y when guard x y -> observe 3
  | K2 _ -> observe 4

```

We note that the right hand side of *observe* is just an arbitrary value and in this case just enumerates the order in which expressions appear. Even if this is an oversimplification of the problem for the prototype, it must be noted that at the compiler level we have the possibility to compile the pattern clauses in two separate steps so that the guards and right-hand-side expressions are semantically equal to their counterparts at the target program level.

```

let _ = function
  | K1 -> lambda0
  | K2 K1 -> lambda1
  | K2 x when lambda-guard0 -> lambda2
  | K2 (K2 x) as y when lambda-guard1 -> lambda3
  | K2 _ -> lambda4

```

### 3.1.1 Matrix decomposition of pattern clauses

We define a new object, the *clause matrix*  $P \rightarrow L$  of size  $|m \times n+1|$  that associates pattern vectors  $\vec{p}_i$  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 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 the *clause matrix*  $P \rightarrow L$ . 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 an empty sequence and 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:

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

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

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 ∈ Set of constructors
  for i ← 1 .. m
    let ki ← constructor_of(pi,1)
    if ki = c then
      P ← qi,1, qi,2, ..., qi,n', pi,2, pi,3, ..., pi,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

```

3. 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

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

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

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

4. 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.

It is important to note that the application of the decomposition algorithm converges. This intuition can be verified by defining the size of the clause matrix  $P \rightarrow L$  as equal to the length of the longest pattern vector  $\vec{p}_i$  and the length of a pattern vector as the number of symbols that appear in the clause. While it is very easy to see that the application of rules 1) and 4) produces new matrices of length equal or smaller than the original clause matrix, we can show that:

- with the application of the constructor rule the pattern vector  $\vec{p}_i$  loses one symbol after its decomposition:

$$\begin{aligned} |(p_{i,1} (q_1, q_2, \dots, q_{n'}), p_{i,2}, p_{i,3}, \dots, p_{i,n})| &= n + n' \\ |(q_{i,1}, q_{i,2}, \dots, q_{i,n'}, p_{i,2}, p_{i,3}, \dots, p_{i,n})| &= n + n' - 1 \end{aligned}$$

- with the application of the orpat rule, we add one row to the clause matrix  $P \rightarrow L$  but the length of a row containing an Or-pattern decreases

$$\begin{aligned} |P \rightarrow L| &= \left| \begin{pmatrix} (p_{1,1}|q_{1,1}) & p_{1,2} & \cdots & p_{1,n} \rightarrow l_1 \\ \vdots & \vdots & \ddots & \vdots \rightarrow \vdots \end{pmatrix} \right| = n + 1 \\ |P' \rightarrow L'| &= \left| \begin{pmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n} \rightarrow l_1 \\ q_{1,1} & p_{1,2} & \cdots & p_{1,n} \rightarrow l_1 \\ \vdots & \vdots & \ddots & \vdots \rightarrow \vdots \end{pmatrix} \right| = n \end{aligned}$$

In our prototype we make use of accessors to encode stored values.

```

let value = 1 :: 2 :: 3 :: []
(* that can also be written *)
let value = []
    |> List.cons 3
    |> List.cons 2
    |> List.cons 1

```

$(\text{field } 0 \text{ } x) = 1$   
 $(\text{field } 0 \text{ (field } 1 \text{ } x)) = 2$   
 $(\text{field } 0 \text{ (field } 1 \text{ (field } 1 \text{ } x))} = 3$   
 $(\text{field } 0 \text{ (field } 1 \text{ (field } 1 \text{ (field } 1 \text{ } x))} =$

An *accessor*  $a$  represents the access path to a value that can be reached by deconstructing the scrutinee.

$a ::= \text{Here} \mid \text{n.a}$

The above example, in encoded form:

```

Here = 1
1.Here = 2
1.1.Here = 3
1.1.1.Here = []

```

In our prototype the source matrix  $m_S$  is defined as follows

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

Source matrices are used to build source decision trees  $C_S$ . 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. %%% Give example of thing

$$C_S ::= \text{Leaf } \text{bb} \mid \text{Switch}(a, (K_i \rightarrow C_i)^{i \in S}, C?) \mid \text{Failure} \mid \text{Unreachable}$$

$$v_S ::= K(v_i)^{i \in I} \mid n \in \mathbb{N}$$

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

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

We define the decision tree of source programs  $\llbracket t_S \rrbracket$  in terms of the decision tree of pattern matrices  $\llbracket m_S \rrbracket$  by the following:

$$\llbracket (p_i \rightarrow bb_i)^{i \in I} \rrbracket := \llbracket (\text{Here}), (p_i \rightarrow bb_i)^{i \in I} \rrbracket$$

Decision tree computed from pattern matrices respect the following invariant:

$$\begin{aligned} \forall v (v_i)^{i \in I} = v(a_i)^{i \in I} &\rightarrow \llbracket m \rrbracket(v) = m(v_i)^{i \in I} \text{ for } m = ((a_i)^{i \in I}, (r_i)^{i \in I}) \\ v(\text{Here}) &= v \\ K(v_i)^i(\text{k.a}) &= v_k(a) \text{ if } k \in [0; n[ \end{aligned}$$

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

Base cases:

1.  $\llbracket \emptyset, (\emptyset \rightarrow bb_i)^i \rrbracket \equiv \text{Leaf } bb_i$  where  $i := \min(I)$ , that is a decision tree  $\llbracket ms \rrbracket$  defined by an empty accessor and empty patterns pointing to blackboxes  $bb_i$ . This respects the invariant because a source matrix in the case of empty rows returns the first expression and  $(\text{Leaf } bb)(v) := \text{Match } bb$
2.  $\llbracket (a_j)^j, \emptyset \rrbracket \equiv \text{Failure}$

Regarding non base cases: Let's first define some auxiliary functions

- The index family of a constructor

$$\text{Idx}(K) := [0; \text{arity}(K)[$$

- head of an ordered family (we write  $x$  for any object here, value, pattern etc.)

$$\text{head}((x_i)^{i \in I}) = x_{\min(I)}$$

- tail of an ordered family

$$\text{tail}((x_i)^{i \in I}) := (x_i)^{i \neq \min(I)}$$



- head constructor of a value or pattern

$$\begin{aligned}\text{constr}(K(x_i)^i) &= K \\ \text{constr}(\_) &= \perp \\ \text{constr}(x) &= \perp\end{aligned}$$

- first non- $\perp$  element of an ordered family

$$\begin{aligned}\text{First}((x_i)^i) &:= \perp && \text{if } \forall i, x_i = \perp \\ \text{First}((x_i)^i) &:= x_{\text{min}}\{i \mid x_i \neq \perp\} && \text{if } \exists i, x_i \neq \perp\end{aligned}$$

- definition of group decomposition:

$$\begin{aligned}\text{let } \text{constrs}((p_i)^{i \in I}) &= \{ K \mid K = \text{constr}(p_i), i \in I \} \\ \text{let } \text{Groups}(m) &\text{ where } m = ((a_i)^i ((p_{ij})^i \rightarrow e_j)^{ij}) = \\ &\text{let } (K_k)^k = \text{constrs}(p_{i0})^i \text{ in} \\ &(\text{ } K_k \rightarrow \\ &\quad ((a_{0..l})^l + \text{tail}(a_i)^i) \\ &\quad ( \\ &\quad \text{if } p_{0j} \text{ is } K_k(q_l) \text{ then} \\ &\quad \quad (q_l)^l + \text{tail}(p_{ij})^i \rightarrow e_j \\ &\quad \text{if } p_{0j} \text{ is } \_ \text{ then} \\ &\quad \quad (\_)^l + \text{tail}(p_{ij})^i \rightarrow e_j \\ &\quad \text{else } \perp \\ &\quad )^j \\ &), ( \\ &\quad \text{tail}(a_i)^i, (\text{tail}(p_{ij})^i \rightarrow e_j \text{ if } p_{0j} \text{ is } \_ \text{ else } \perp)^j \\ &)\end{aligned}$$

$\text{Groups}(m)$  is an auxiliary function that decomposes a matrix  $m$  into submatrices, according to the head constructor of their first pattern.  $\text{Groups}(m)$  returns one submatrix  $m\_r$  for each head constructor  $K$  that occurs on the first row of  $m$ , plus one "wildcard submatrix"  $m_{\text{wild}}$  that matches on all values that do not start with one of those head constructors. Intuitively,  $m$  is equivalent to its decomposition in the following sense: if the first pattern of

an input vector  $(v_i)^i$  starts with one of the head constructors  $K_k$ , then running  $(v_i)^i$  against  $m$  is the same as running it against the submatrix  $m_{K_k}$ ; otherwise (its head constructor  $\notin (K_k)^k$ ) it is equivalent to running it against the wildcard submatrix.

We formalize this intuition as follows

1. Lemma (Groups): Let  $m$  be a matrix with

$$\text{Groups}(m) = (k_r \rightarrow m_r)^k, m_{\text{wild}}$$

For any value vector  $(v_i)^i$  such that  $v_0 = k(v'_i)^i$  for some constructor  $k$ , we have:

$$\begin{aligned} &\text{if } k = k_k \text{ for some } k \text{ then} \\ &\quad m(v_i)^i = m_k((v'_i)^i + (v_i)^{i \in \Gamma \setminus \{0\}}) \\ &\text{else} \\ &\quad m(v_i)^i = m_{\text{wild}}(v_i)^{i \in \Gamma \setminus \{0\}} \end{aligned}$$

2. Proof: Let  $m$  be a matrix  $((a_i)^i, ((p_{ij})^i \rightarrow e_j)^j)$  with

$$\text{Groups}(m) = (K_k \rightarrow m_k)^k, m_{\text{wild}}$$

Below we are going to assume that  $m$  is a simplified matrix such that the first row does not contain an or-pattern or a binding to a variable.

Let  $(v_i)^i$  be an input matrix with  $v_0 = K_v(v'_i)^i$  for some constructor  $K_v$ . We have to show that:

- if  $K_k = K_v$  for some  $K_k \in \text{constrs}(p_{0j})^j$ , then

$$m(v_i)^i = m_k((v'_i)^i + \text{tail}(v_i)^i)$$

- otherwise  $m(v_i)^i = m_{\text{wild}}(\text{tail}(v_i)^i)$

Let us call  $(r_{kj})$  the  $j$ -th row of the submatrix  $m_k$ , and  $r_{j\text{wild}}$  the  $j$ -th row of the wildcard submatrix  $m_{\text{wild}}$ .

Our goal contains same-behavior equalities between matrices, for a fixed input vector  $(v_i)^i$ . It suffices to show same-behavior equalities between each row of the matrices for this input vector. We show that for any  $j$ ,

- if  $K_k = K_v$  for some  $K_k \in \text{constrs}(p_{0j})^j$ , then

$$(p_{ij})^i(v_i)^i = r_{kj}((v'_l)^l + \text{tail}(v_i)^i)$$

- otherwise

$$(p_{ij})^i(v_i)^i = r_{j\text{wild}} \text{tail}(v_i)^i$$

In the first case ( $K_v$  is  $K_k$  for some  $K_k \in \text{constrs}(p_{0j})^j$ ), we have to prove that

$$(p_{ij})^i(v_i)^i = r_{kj}((v'_l)^l + \text{tail}(v_i)^i)$$

By definition of  $m_k$  we know that  $r_{kj}$  is equal to

$$\begin{aligned} &\text{if } p_{0j} \text{ is } K_k(q_l) \text{ then} \\ &\quad (q_l)^l + \text{tail}(p_{ij})^i \rightarrow e_j \\ &\text{if } p_{0j} \text{ is } \_ \text{ then} \\ &\quad (\_)^l + \text{tail}(p_{ij})^i \rightarrow e_j \\ &\text{else } \perp \end{aligned}$$

By definition of  $(p_i)^i(v_i)^i$  we know that  $(p_{ij})^i(v_i)^i$  is equal to

$$\begin{aligned} &(K(q_j)^j, \text{tail}(p_{ij})^i) (K(v'_l)^l, \text{tail}(v_i)^i) := ((q_j)^j + \text{tail}(p_{ij})^i)((v'_l)^l + \text{tail}(v_i)^i) \\ &(\_, \text{tail}(p_{ij})^i) (v_i) := \text{tail}(p_{ij})^i(\text{tail}(v_i)^i) \\ &(K(q_j)^j, \text{tail}(p_{ij})^i) (K'(v'_l)^l, \text{tail}(v_j)^j) := \perp \text{ if } K \neq K' \end{aligned}$$

We prove this first case by a second case analysis on  $p_{0j}$ .

TODO

In the second case ( $K_v$  is distinct from  $K_k$  for all  $K_k \in \text{constrs}(p_{0j})^j$ ), we have to prove that

$$(p_{ij})^i (v_i)^i = r_{j\text{wild}} \text{tail}(v_i)^i$$

TODO

### 3.2 Target translation

The target 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 a subset of the Lambda intermediate language, down to OCaml code. This is the grammar of the target language (TODO: check menhir grammar)

```

start ::= sexpr
sexpr ::= variable | string | "(" special_form ")"
string ::= "\" identifier \"" ;; string between doublequotes
variable ::= identifier
special_form ::= let|catch|if|switch|switch-star|field|apply|isout
let ::= "let" assignment sexpr ;; (assignment sexpr)+ outside of pattern match code
assignment ::= "function" variable variable+ ;; the first variable is the identifier of the function
field ::= "field" digit variable
apply ::= ocaml_lambda_code ;; arbitrary code
catch ::= "catch" sexpr with sexpr
with ::= "with" "(" label ")"
exit ::= "exit" label
switch-star ::= "switch*" variable case*
switch ::= "switch" variable case* "default:" sexpr
case ::= "case" casevar ":" sexpr
casevar ::= ("tag"|int) integer
if ::= "if" bexpr sexpr sexpr
bexpr ::= "(" ("!="|"="|\vert{ }">"|<="|">"|<") sexpr digit | field ")"
label ::= integer

```

The prototype doesn't support strings.

The AST built by the parser is traversed and evaluated by the symbolic execution engine. Given that the target language supports jumps in the form of "catch - exit" blocks the engine tries to evaluate the instructions inside the blocks and stores the result of the partial evaluation into a record. When a jump is encountered, the information at the point allows to finalize the evaluation of the jump block. In the environment the engine also stores bindings to values and functions. Integer additions and subtractions are simplified in a second pass. The result of the symbolic evaluation is a target decision tree  $C_T$

$$C_T ::= \text{Leaf } bb \mid \text{Switch}(a, (\pi_i \rightarrow C_i)^{i \in S}, C?) \mid \text{Failure}$$

$$v_T ::= \text{Cell}(\text{tag} \in \mathbb{N}, (v_i)^{i \in I}) \mid n \in \mathbb{N}$$

Every branch of the decision tree is "constrained" by a domain

$$\text{Domain } \pi = \{ n \mid n \in \mathbb{N} \times n \mid n \in \text{Tag} \subseteq \mathbb{N} \}$$

Intuitively, the  $\pi$  condition at every branch tells us the set of possible values that can "flow" through that path.  $\pi$  conditions are refined by the engine during the evaluation; at the root of the decision tree the domain corresponds to the set of possible values that the type of the function can hold.  $C?$  is the fallback node of the tree that is taken whenever the value at that point of the execution can't flow to any other subbranch. Intuitively, the  $\pi_{\text{fallback}}$  condition of the  $C?$  fallback node is

$$\pi_{\text{fallback}} = \neg \bigcup_{i \in I} \pi_i$$

The fallback node can be omitted in the case where the domain of the children nodes correspond to set of possible values pointed by the accessor at that point of the execution

$$\text{domain}(v_S(a)) = \bigcup_{i \in I} \pi_i$$

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

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

### 3.3 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. Our algorithm respects the following correctness statement:

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

Our equivalence-checking algorithm  $\text{equiv}(S, C_S, C_T)G$  is an exactly decision procedure for the provability of the judgment  $(\text{equiv}(S, C_S, C_T)G)$ , defined

below.

$$\begin{array}{ll}
\text{constraint trees} & \text{boolean result} \\
C ::= \text{Leaf}(t) & b \in \{0, 1\} \\
| \text{Failure} & \\
| \text{Switch}(a, (\pi_i, C_i)^i, C_{\text{fb}}) & \text{guard queues} \\
| \text{Guard}(t, C_0, C_1) & G ::= (t_1 = b_1), \dots, (t_n = b_n)
\end{array}$$

$$\begin{array}{c}
\text{input space} \\
S \subseteq \{(v_S, v_T) \mid v_S \approx_{\text{val}} v_T\} \\
\hline
\text{equiv}(\emptyset, C_S, C_T)G \qquad \text{equiv}(S, \text{Failure}, \text{Failure})\square
\end{array}$$

$$\frac{t_S \approx_{\text{term}} t_T}{\text{equiv}(S, \text{Leaf}(t_S), \text{Leaf}(t_T))\square}$$

$$\frac{\forall i, \text{equiv}((S \wedge a = K_i), C_i, \text{trim}(C_T, a = K_i))G \quad \text{equiv}((S \wedge a \notin (K_i)^i), C_{\text{fb}}, \text{trim}(C_T, a \notin (K_i)^i))G}{\text{equiv}(S, \text{Switch}(a, (K_i, C_i)^i, C_{\text{fb}}), C_T)G}$$

$$\frac{C_S \in \text{Leaf}(t), \text{Failure} \quad \forall i, \text{equiv}((S \wedge a \in D_i), C_S, C_i)G \quad \text{equiv}((S \wedge a \notin (D_i)^i), C_S, C_{\text{fb}})G}{\text{equiv}(S, C_S, \text{Switch}(a, (D_i)^i C_i, C_{\text{fb}}))G}$$

$$\frac{\text{equiv}(S, C_0, C_T)G, (t_S = 0) \quad \text{equiv}(S, C_1, C_T)G, (t_S = 1)}{\text{equiv}(S, \text{Guard}(t_S, C_0, C_1), C_T)G}$$

$$\frac{t_S \approx_{\text{term}} t_T \quad \text{equiv}(S, C_S, C_b)G}{\text{equiv}(S, C_S, \text{Guard}(t_T, C_0, C_1))(t_S = b), G}$$

Running a program  $t_S$  or its translation  $\llbracket t_S \rrbracket$  against an input  $v_S$  produces as a result  $r$  in the following way:

$$\begin{array}{l}
(\llbracket t_S \rrbracket_S(v_S) \equiv C_S(v_S)) \rightarrow r \\
t_S(v_S) \rightarrow r
\end{array}$$

Likewise

$( \llbracket t_T \rrbracket_T(v_T) \equiv C_T(v_T) ) \rightarrow r$   
 $t_T(v_T) \rightarrow r$   
 result  $r ::=$  guard list \* (Match blackbox | NoMatch | Absurd)  
 guard  $::=$  blackbox.

Having defined equivalence between two inputs of which one is expressed in the source language and the other in the target language,  $v_S \simeq v_T$ , 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 result of the proposed equivalence algorithm is *Yes* or *No*( $v_S, v_T$ ). In particular, in the negative case,  $v_S$  and  $v_T$  are a couple of possible counter examples for which the decision trees produce a different result.

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:

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

We say that  $C_T$  covers the input space  $S$ , written  $covers(C_T, S)$  when every value  $v_S \in S$  is a valid input to the decision tree  $C_T$ . (TODO: rephrase) Given an input space  $S$  and a couple of decision trees, where the target decision tree  $C_T$  covers the input space  $S$  we can define equivalence:

$$equiv(S, C_S, C_T, gs) = Yes \wedge covers(C_T, S) \rightarrow \forall v_S \simeq v_T \in S, C_S(v_S) \simeq_{gs} C_T(v_T)$$

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

$$equiv(S, C_S, C_T, gs) = No(v_S, v_T) \wedge covers(C_T, S) \rightarrow v_S \simeq v_T \in S \wedge C_S(v_S) \neq_{gs} C_T(v_T)$$

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

$$equiv(S, \llbracket t_S \rrbracket_S, \llbracket t_T \rrbracket_T, \emptyset) = Yes \Leftrightarrow t_S \simeq t_T$$



### 3.3.1 The trimming lemma

The trimming lemma allows to reduce the size of a decision tree given an accessor  $a \rightarrow \pi$  relation (TODO: expand)

$$\forall v_T \in (a \rightarrow \pi), C_T(v_T) = C_{t/a \rightarrow \pi}(v_T)$$

We prove this by induction on  $C_T$ :

- $C_T = \text{Leaf}_{\text{bb}}$ : when the decision tree is a leaf terminal, the result of trimming on a Leaf is the Leaf itself

$$\text{Leaf}_{\text{bb}/a \rightarrow \pi}(v) = \text{Leaf}_{\text{bb}}(v)$$

- The same applies to Failure terminal

$$\text{Failure}_{/a \rightarrow \pi}(v) = \text{Failure}(v)$$

- When  $C_T = \text{Switch}(b, (\pi \rightarrow C_i)^i)_{/a \rightarrow \pi}$  then we look at the accessor  $a$  of the subtree  $C_i$  and we define  $\pi_i' = \pi_i$  if  $a \neq b$  else  $\pi_i \cap \pi$ . Trimming a switch node yields the following result:

$$\text{Switch}(b, (\pi \rightarrow C_i)^{i \in I})_{/a \rightarrow \pi} := \text{Switch}(b, (\pi_i' \rightarrow C_{i/a \rightarrow \pi})^{i \in I})$$

For the trimming lemma we have to prove that running the value  $v_T$  against the decision tree  $C_T$  is the same as running  $v_T$  against the tree  $C_{\text{trim}}$  that is the result of the trimming operation on  $C_T$

$$C_T(v_T) = C_{\text{trim}}(v_T) = \text{Switch}(b, (\pi_i' \rightarrow C_{i/a \rightarrow \pi})^{i \in I})(v_T)$$

We can reason by first noting that when  $v_T \notin (b \rightarrow \pi_i)^i$  the node must be a Failure node. In the case where  $\exists k \mid v_T \in (b \rightarrow \pi_k)$  then we can prove that

$$C_{k/a \rightarrow \pi}(v_T) = \text{Switch}(b, (\pi_i' \rightarrow C_{i/a \rightarrow \pi})^{i \in I})(v_T)$$

because when  $a \neq b$  then  $\pi_k' = \pi_k$  and this means that  $v_T \in \pi_k'$  while when  $a = b$  then  $\pi_k' = (\pi_k \cap \pi)$  and  $v_T \in \pi_k'$  because:

- by the hypothesis,  $v_T \in \pi$
- we are in the case where  $v_T \in \pi_k$

So  $v_T \in \pi_k$  and by induction

$$C_k(v_T) = C_{k/a \rightarrow \pi}(v_T)$$

We also know that  $\forall v_T \in (b \rightarrow \pi_k) \rightarrow C_T(v_T) = C_k(v_T)$  By putting together the last two steps, we have proven the trimming lemma.

### 3.3.2 Equivalence checking

The equivalence checking algorithm takes as parameters an input space  $S$ , a source decision tree  $C_S$  and a target decision tree  $C_T$ :

$$\text{equiv}(S, C_S, C_T) \rightarrow \text{Yes} \mid \text{No}(v_S, v_T)$$

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:

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

where  $\mathbf{k}$  is a constructor.

We proceed by case analysis:

1. in case of unreachable:

$$C_S(v_S) = \text{Absurd}(\text{Unreachable}) \neq C_T(v_T) \ \forall v_S, v_T$$

1. In the case of an empty input space

$$\text{equiv}(\emptyset, C_S, C_T) := \text{Yes}$$

and that is trivial to prove because there is no pair of values  $(v_S, v_T)$  that could be tested against the decision trees. In the other subcases  $S$  is always non-empty.

2. When there are *Failure* nodes at both sides the result is *Yes*:

$$\text{equiv}(S, \text{Failure}, \text{Failure}) := \text{Yes}$$

Given that  $\forall v, \text{Failure}(v) = \text{Failure}$ , the statement holds.

3. When we have a Leaf or a Failure at the left side:

$$\begin{aligned}
&\text{equiv}(S, \text{Failure as } C_S, \text{Switch}(a, (\pi_i \rightarrow C_{T_i})^{i \in I})) := \text{Forall}(\text{equiv}(S \cap a \rightarrow \pi(k_i)), C_S, C_{T_i})^{i \in I} \\
&\text{equiv}(S, \text{Leaf } bb_S \text{ as } C_S, \text{Switch}(a, (\pi_i \rightarrow C_{T_i})^{i \in I})) := \text{Forall}(\text{equiv}(S \cap a \rightarrow \pi(k_i)), C_S, C_{T_i})^{i \in I}
\end{aligned}$$

Our algorithm either returns *Yes* for every sub-input space  $S_i := S \cap (a \rightarrow \pi(k_i))$  and subtree  $C_{T_i}$

$$\text{equiv}(S_i, C_S, C_{T_i}) = \text{Yes} \ \forall i$$

or we have a counter example  $v_S, v_T$  for which

$$v_S \simeq v_T \in S_k \wedge c_S(v_S) \neq C_{T_k}(v_T)$$

then because

$$\begin{aligned} v_T \in (a \rightarrow \pi_k) &\rightarrow C_T(v_T) = C_{T_k}(v_T), \\ v_S \simeq v_T \in S &\wedge C_S(v_S) \neq C_T(v_T) \end{aligned}$$

we can say that

$$\text{equiv}(S_i, C_S, C_{T_i}) = \text{No}(v_S, v_T) \text{ for some minimal } k \in I$$

4. When we have a Switch on the right we define  $\pi_{\text{fallback}}$  as the domain of values not covered but the union of the constructors  $k_i$

$$\pi_{\text{fallback}} = \neg \bigcup_{i \in I} \pi(k_i)$$

Our algorithm proceeds by trimming

$$\begin{aligned} &\text{equiv}(S, \text{Switch}(a, (k_i \rightarrow C_{S_i})^{i \in I}, C_{\text{sf}}), C_T) := \\ &\text{Forall}(\text{equiv}(S \cap (a \rightarrow \pi(k_i)^{i \in I}), C_{S_i}, C_{t/a \rightarrow \pi(k_i)})^{i \in I} + \text{equiv}(S \cap (a \rightarrow \pi_n), C_S, C_{a \rightarrow \pi_{\text{fallback}}})) \end{aligned}$$

The statement still holds and we show this by first analyzing the *Yes* case:

$$\text{Forall}(\text{equiv}(S \cap (a \rightarrow \pi(k_i)^{i \in I}), C_{S_i}, C_{t/a \rightarrow \pi(k_i)})^{i \in I} = \text{Yes}$$

The constructor  $k$  is either included in the set of constructors  $k_i$ :

$$k \mid k \in (k_i)^i \wedge C_S(v_S) = C_{S_i}(v_S)$$

We also know that

$$\begin{aligned} (1) \quad &C_{S_i}(v_S) = C_{t/a \rightarrow \pi_i}(v_T) \\ (2) \quad &C_{T/a \rightarrow \pi_i}(v_T) = C_T(v_T) \end{aligned}$$

(1) is true by induction and (2) is a consequence of the trimming lemma. Putting everything together:

$$C_S(v_S) = C_{S_i}(v_S) = C_{T/a \rightarrow \pi_i}(v_T) = C_T(v_T)$$

When the  $k \notin (k_i)^i$  [TODO]

The auxiliary Forall function returns  $No(v_S, v_T)$  when, for a minimum  $k$ ,

$$\text{equiv}(S_k, C_{S_k}, C_{T/a \rightarrow \pi_k} = No(v_S, v_T)$$

Then we can say that

$$C_{S_k}(v_S) \neq C_{t/a \rightarrow \pi_k}(v_T)$$

that is enough for proving that

$$C_{S_k}(v_S) \neq (C_{t/a \rightarrow \pi_k}(v_T) = C_T(v_T))$$