



**HAL**  
open science

# Knit&Frog: Pattern matching compilation for custom memory representations

Thaïs Baudon, Laure Gonnord, Gabriel Radanne

► **To cite this version:**

Thaïs Baudon, Laure Gonnord, Gabriel Radanne. Knit&Frog: Pattern matching compilation for custom memory representations. [Research Report] Inria Lyon. 2022. hal-03684334v1

**HAL Id: hal-03684334**

**<https://inria.hal.science/hal-03684334v1>**

Submitted on 1 Jun 2022 (v1), last revised 20 Jul 2022 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Knit&Frog: Pattern matching compilation for custom memory representations

Thaïs BAUDON, Laure GONNORD, Gabriel RADANNE

**RESEARCH  
REPORT**

**N° 9473**

June 2022

Project-Team CASH

ISSN INRIA/RR--9473--FR+ENG

ISSN 0249-6399





## Knit&Frog: Pattern matching compilation for custom memory representations

Thaïs BAUDON\*, Laure GONNORD†, Gabriel RADANNE‡

Project-Team CASH

Research Report n° 9473 — June 2022 — 36 pages

**Abstract:** Initially present only in functional languages such as OCaml and Haskell, Algebraic Data Types have now become pervasive in mainstream languages, providing nice data abstractions and an elegant way to express functions through *pattern-matching*. Numerous approaches have been designed to compile rich pattern matching to cleverly designed, efficient decision trees. However, these approaches are specific to a choice of *internal memory representation* which must accommodate garbage-collection and polymorphism.

ADTs now appear in languages more liberal in their memory representation such as Rust. Notably, Rust is now introducing more and more optimizations of the memory layout of Algebraic Data Types. As memory representation and compilation are interdependent, it raises the question of pattern matching compilation in the presence of non-regular, potentially customized, memory layouts.

In this report, we present Knit&Frog, a framework to compile pattern-matching for monomorphic ADTs, *parametrized* by an arbitrary memory representation. We propose a novel way to describe choices of memory representation along with a validity condition under which we prove the correctness of our compilation scheme. The approach is implemented in a prototype tool *ribbit*.

**Key-words:** Algebraic Data Types, Pattern Matching, Compilation

---

\* ENS de Lyon, LIP (UMR CNRS/ENS Lyon/UCB Lyon1/INRIA),F-69000 Lyon, France

† Grenoble-INP/LCIS & LIP

‡ Inria, LIP (UMR CNRS/ENS Lyon/UCB Lyon1/INRIA),F-69000 Lyon, France

## Knit&Frog: Pattern matching compilation for custom memory representations

**Résumé :** Ce rapport de recherche présente Knit&Frog, un cadre formel pour décrire des représentations mémoires de types algébriques, ainsi que son usage dans la formalisation d'un algorithme de filtrage (*pattern-matching*) paramétré par la représentation mémoire. La correction de l'algorithme est prouvé sous une hypothèse de *validité* de la représentation. L'algorithme est implémenté dans l'outil `ribbit`.

**Mots-clés :** HPC, types algébriques, filtrage, compilation

# 1 Introduction

Algebraic Data Types (ADTs) are an essential tool to model data and information. They allow to group together information in a consistent way through the use of records, also called product types, and to organize options through the use of variants, also called sum types. Product and sum types together allow to enforce invariants present in the domain under study and provide convenient tools to inspect these data-types through pattern matching.

Combined, these features offers numerous advantages:

- Model data in a way that is close to the programmer’s intuition, abstracting away the details of the memory representation of said data.
- Safely handle data by ensuring via pattern-matching that its manipulation is well-typed, exhaustive and non-redundant.
- Optimize manipulation of data thanks to rich constructs understood by the compiler.

Despite these promises, Algebraic Data Types were initially only present in functional programming languages such as OCaml and Haskell. Recently, they have gained a foothold in more mainstream languages such as Typescript, Scala, Rust and even soon Java. They are however still lacking in high-performance lower level languages. One difficulty to language designers who want to add pattern matching to their language is that compiling a rich pattern matching language to efficient code is a non-trivial task, which is not commonly available in shared compiler frameworks such as LLVM. Indeed, such frameworks only provide optimizations for C-like switches on integers (or integer-like enumerations). Additionally, existing works on pattern matching Maranget (2008); Wadler (1987); Sestoft (1996) provide very efficient compilation schemes, but are geared towards memory representations found in GC-managed functional languages such as OCaml and Haskell: uniform representations with liberal usage of boxing. Highly non-uniform data representations such as the ones found in C++ do not easily fit.

More generally, the descriptive nature of ADTs should enable compilers to aggressively optimize the representation of terms. The simplest example is the `Option` type, which is either `Some value` or `None`. If the value in question is an integer from 0 to 10, `None` can easily be represented as 11. This optimization is regularly done by programmers manually, forgoing the guarantees provided by languages. More complex optimizations on nested and rich data-types are even more error-prone. These transformations could easily be done automatically by the compiler. This trove of optimization potential has been brushed on in recent versions of Rust, but remains largely unexplored.

Finally, on top of the previously mentioned difficulty of adapting pattern matching compilation schemes to irregular memory representations, the correctness of such compilation schemes is also delicate to establish when the terms can be arbitrarily shaped.

In this article, we lay the groundwork for highly optimized pattern matching in any language with arbitrary non-uniform memory representation. We present Knit&Frog, a compilation framework for rich pattern languages with arbitrary memory representation of monomorphic terms:

- We provide a characterization of memory representations for terms of Algebraic Data Types as two main operations: **Knit** which builds a term into its memory representation, and **Frog** which deconstructs the memory representation to identify the underlying term.
- We provide a compilation scheme parameterized by the memory representation. This scheme is adapted from Maranget (2008), the best implementation of pattern matching known thus far.
- We provide a sufficient condition for a memory representation to yield itself to pattern matching and prove end-to-end correctness of our compilation scheme in this case.

We start by giving some examples of memory representations and pattern-matching compilation as motivation in Section 2. We then provide a formal presentation of our pattern language (Section 3) and our target memory representations (Section 4). Finally, we present our compilation scheme (Section 5) along with end-to-end examples (Section 6) and a proof of correctness (Section 7).

```

1 enum Option<T> {
2   None, // No value
3   Some (T), // Some value
4 }
5 type Interval = Option<&(u64, u64)>

1 fn length(x: Interval) -> u64 {
2   match x {
3     None => 0,
4     Some(&(i1, j1)) => j1 - i1
5   }}

```

Figure 1: Example of program manipulating the `Option` type in Rust

## 2 Algebraic Data Types and their memory representation

Let us now explore the different memory representation choices for Algebraic Data Types (ADTs, for short), and how their values are manipulated. In this section, we present various examples of types and programs and discuss how they are currently represented in existing languages, and how they might be. For illustration, we will look at the output of two languages which implement ADTs: Rust and OCaml. While these languages have some common lineage, they have a very different attitude towards code emission: OCaml is a GC-managed language which factors predictability and regularity. Rust on the other hand favors performance and absolute control over low-level details. These differences result in drastically different choices in memory representation. We thus start our exploration by a classic question: how to represent the option type?

**Example 1** (Option type). We consider the program in Fig. 1 in Rust syntax. We first restate the `Option<T>` type, whose goal is to model that some data of type `T` is *optionally* present. Either we have `Some` data, or `None`. This is modeled by the two constructors `None` and `Some`. The `None` constructor is argument-less, while the `Some` constructor has a parameter of type `T`.

The `Interval` type is an optional pair of integers. `u64` denotes an unsigned integer over 64 bits and `&(u64, u64)` denotes a *reference*<sup>1</sup> to a pair of such integers. Finally, `Option<&(u64, u64)>` denotes an optional pair. The value `Some(&(u1756, u1791))` is of type `Interval` and designates an interval from 1756 to 1791.

The `length` function takes an interval and returns its length using *pattern matching*. If its argument is of the “shape” of the left-hand side of the rule then the value of the right-hand side (body) is evaluated. In Rust, pattern matching is introduced by the keyword `match` and alternatives are depicted under the form of a list of the form  $p \Rightarrow b$  where  $p$  is a *pattern* and  $b$  its body. Moreover, patterns can be nested, and the body can use named subterms. In our example, `None` yields a length of 0 and `Some(i, j)` yields a length of  $j - i$ .

We now look at how Rust and OCaml compile the `length` function. Both compilers provide an intermediate representation (IR) in which the pattern matching is represented as a decision tree which manipulates the underlying memory representation. We showcase simplified versions of the parts of interest in Fig. 2. OCaml’s Lambda IR represents matching using the built-in primitive `field` to access subterms inside the representation. It uses the `if` primitive to discriminate between cases. Rust’s MIR<sup>2</sup> uses slightly lower level operations such as dereferencing and field accesses, and a `switch` construct on memory words.

In both cases, `None` is represented as an unboxed integer. The similarity stops there. In OCaml, sum types are represented uniformly: argument-less constructors are unboxed integers, and constructors with arguments are pointers to a block, which always starts with a tag and some fields. This block contains a single field with a pointer to the inner tuple. This tuple is a block with a tag and two fields, containing the integers. This totals 6 memory words<sup>3</sup>. In Rust, `Some` is a pointer referencing directly a memory span of 2 memory words. Rust recognizes `Option<&T>` as something that can be represented compactly, by remarking that Rust pointers are non-null. It can thus use 0 as the `None` value, and directly use the pointer in the `Some` cases. The resulting type uses 3 memory words.

Our goal here is not to oppose the two representations: the OCaml representation is designed to support metadata for the GC and works smoothly with polymorphism and separate compilation, while Rust’s representation focuses on efficiency and control. Both are excellent designs for their objectives. However, in both cases, we want efficient implementations of pattern matching to deconstruct values in memory. In fact, the choice of memory representation can have a great influence on the shape of the code generated to match values (in terms of number of switches, complexity of expressions...). The goal of this article is precisely to

<sup>1</sup>For simplicity of presentation, we omit details regarding lifetimes and differences between `&` and `Box`.

<sup>2</sup>the Middle Intermediate Representation is a control flow graph. We reconstruct some of its logic here for ease of reading.

<sup>3</sup>Knowledgeable readers might remark that it is possibly to give up the common `Option` type and inline the definition of the tuple inside the sum, going down to 1 indirection and 4 memory words at the price of modularity.

<pre> 1 fn length(_1:Interval) = 2   switch(_1){ 3     0 -&gt; 0 4     1 -&gt; let _5 = *((_1 as Some).0).0; 5           let _6 = *((_1 as Some).0).1; 6           _6 - _5 7 }                 </pre> <p>(a) Simplified Rust “MIR” output</p>	<pre> 1 length(param/85) = 2   (if (== param/85 1) 3     0 4     (let (*match*/88 (field 1 param/85)) 5       (- (field 2 *match*/88) (field 1 *match*/88)) 6     )) 7 }                 </pre> <p>(b) Simplified OCaml “Lambda” output</p>
---	---

Figure 2: Intermediate representation of the length function in Rust and OCaml

<pre> 1 //Colors for the Red-Black Trees 2 enum Color { Red, Black, } 3 4 //Red-Black Trees of content T 5 enum RBT&lt;T&gt; { 6   Node (Color, T, &amp;RBT&lt;T&gt;, &amp;RBT&lt;T&gt;), 7   Empty, 8 }                 </pre>	<pre> 1 match c, v, t1, t2 { 2   Black, z, &amp;Node(Red, y, &amp;Node(Red, x, a, b), c), d 3     Black, z, &amp;Node(Red, x, a, &amp;Node(Red, y, b, c)), d 4     Black, x, a, &amp;Node(Red, z, &amp;Node(Red, y, b, c), d) 5     Black, x, a, &amp;Node(Red, y, b, &amp;Node(Red, z, c, d)) 6   =&gt; Node(Red, y, &amp;Node(Black, x, a, b), &amp;Node(Black, z, c, d)), 7   a, b, c, d =&gt; Node (a, b, c, d), 8 }                 </pre>
---	---

Figure 3: Example of Red-Black tree and their balancing operation in Rust

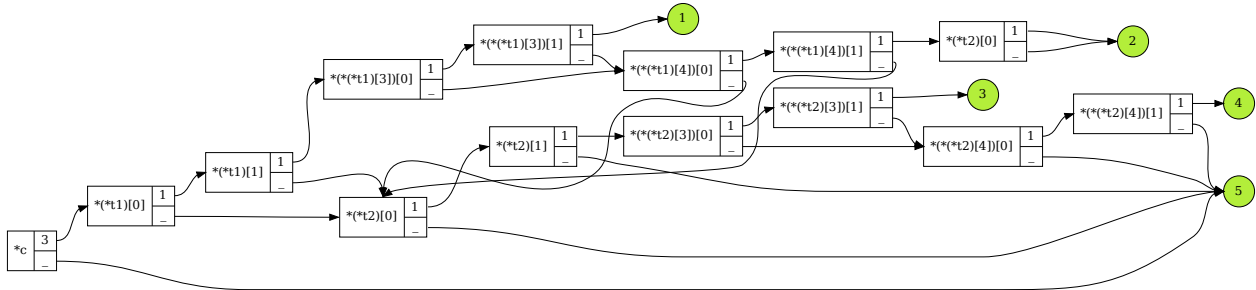


Figure 4: Output of our ribbit compiler for the code of Fig. 3 in the OCaml memory representation

This output is a (decision tree) whose root is on the left. Decision (square) nodes compute values from their input (low-level representation), then branch depending on this value. Round nodes denote the final returned value of the procedure (the numbering of the matched branch).

provide efficient compilation of complex patterns for any memory representation. As patterns become more complex, the compiler has to do more work, as we show in Example 2.

**Example 2** (Search Trees). To illustrate the need for a sophisticated pattern matching compiler, we now look at *Red-Black Trees*. A Rust version is depicted in Fig. 3. This type definition expresses trees as recursive data structures, and uses a tuple for the various fields in the `Node` case. For instance, `Node(Red, 1515, &Node(Black, &Empty,&Empty), &Empty)` is a tree of type `RBT<Int>`.

Red-Black trees famously rely on a fairly complex balancing step, which redistributes the colors depending on the internal invariant of the data structure. Thanks to nested pattern matching, this step can be expressed very compactly, as shown in Fig. 3. This pattern matching inspects four arguments at the same time: the current color `c`, the current value `v` and the sub-trees `t1` and `t2`.

Since this pattern matching is at the core of a performance-sensitive data structure, we naturally want it to be as efficient as possible. This is why many pattern matching implementations come with clever heuristics and techniques to output optimized decision trees Kosarev et al. (2020); Maranget (2008); Sestoft (1996). The resulting code is highly non-trivial, as can be seen in Fig. 4.

Such optimized decision trees should still respect, and even take advantage of, the memory representation. For instance, the type `RBT<u64>` is represented by the Rust compiler using only 4 memory words per `Node`. A node is then made of a word with a bit marking it as non-0 (to avoid confusion with `None`) and the color bit, followed by the 64-bit integer and the two pointers. It is however possible to be even more compact: the Linux kernel uses a hand-crafted representation of Red-Black trees that exploits bit-stealing to use one less word Wilke (2016): since pointers are word-aligned, colors can be stored in the lower bits of each pointer. These complex manually-tuned optimizations are common in low-level C programming, but come at a high cost: programmers forgo the use of modern safer constructs such as Algebraic Data Types and pattern



matching. Furthermore, they now need to design the decision tree themselves, and make sure their choice of representation contains enough information to distinguish, for instance, between the `Empty` and `Node` cases. Moreover, the emission of code for such mangled objects in memory is delicate (notably, the bit-stealing technique in the Linux kernel is undefined behavior in the C standard).

Ideally, we would like to still use high-level pattern languages, while the compilation process adapts and exploits the optimized representation. This article makes a first step towards this goal by presenting Knit&Frog, a pattern matching compilation technique **parameterized by the memory representation**. We formalize what it means to define a memory representation, present a validity criterion ensuring its suitability for pattern matching and prove the correctness of our compilation scheme. We also present several examples of both toy and realistic memory representations. Finally, we implemented our technique in a tool called `ribbit`, whose output using the OCaml representation is shown in Fig. 4.

### 3 Pattern language

We now present a formal version of our source language. For clarity and ease of presentation, we make some simplifying assumptions. First off, we only consider the pattern part of the language, which represents the core of our compilation procedure. Furthermore, we assume the following:

- Types are monomorphic. All code has been specialized previously, and all type variables have been replaced by concrete types.
- Patterns are always *exhaustive*: all possible cases are handled.
- Constant types fit in a machine word, such as machine integers, floats, chars,  $\dots$ . Strings are thus not a constant type (and could, for instance, be represented as arrays of chars).
- Types are non-recursive. This limitation is far less restrictive than it seems for compilation purposes since recursive types can be unfolded to the size of the given pattern.

We shall come back to these limitations later in the article.

Our language, described in Fig. 5, roughly follows the syntax of ML-style languages, restricted to simple types and patterns. A matching, denoted  $m$ , is composed of a list of patterns. The bodies of the clauses are left unspecified here, and can be composed of any expression language. Patterns, denoted  $p$ , are composed of constant, product, constructor and reference patterns as introduced previously, along with variables ( $x$ ), wildcards ( $\_$ ) and “or”-patterns ( $p_1 \mid p_2$ ). Types, denoted  $\tau$ , are composed of constant types, product types ( $\prod_k \tau_k$ , also written  $(\tau_1, \tau_2, \dots)$ ), sum types ( $\sum_k K_k(\tau_{k,1}, \dots, \tau_{k,n_k})$  also written  $K_1(\tau_{1,1}, \dots, \tau_{1,n_1}) + K_2(\tau_{2,1}, \dots, \tau_{2,n_2}) + \dots$ ) and reference types ( $\&\tau$ ). Typing environments, denoted  $\Gamma$ , associate variables to their types. Finally, values, denoted  $v$ , follow a subset of the pattern grammar. Value environments, denoted  $\sigma$ , associate variables to their values.

**Example 3** (Interval example, cont’). The `Interval` type of Example 1 is denoted by

$$\text{Interval} = \text{None}() + \text{Some}(\&(u64, u64)),$$

The first constructor `None` has arity 0 and the second one `Some` has arity 1 with a reference type to a tuple (product of size 2). The matching problem for function `length` will be expressed as:

$$\text{Match} \left\{ \begin{array}{l} \mid \text{None} \\ \mid \text{Some}(\&(x_1, x_2)) \end{array} \right\}$$

#### 3.1 Semantics

The semantics of patterns is presented in Fig. 6, which defines the judgment  $p \triangleright v \rightarrow \sigma$  meaning that pattern  $p$  matches value  $v$  and binds environment  $\sigma$ . The object being matched (here,  $v$ ) is also called the *discriminant*. We also assume the existence of the negated judgment,  $p \not\triangleright v$  where pattern  $p$  doesn’t match value  $v$ . Most of the rules are straightforward, with the following specificities:

- The bound value environment returned by the matching is populated by variables, through the VAR rule.

Patterns		Types	
$p ::= c \in \mathcal{C}$	(Constant pattern)	$\tau ::= T$	(Built-in Type)
$\langle p_1, \dots, p_n \rangle$	(Product pattern)	$\prod_n \tau_i$	(Product Type)
$K(p_1, \dots, p_n)$	(Constructor pattern)	$\sum_n K_i(\tau_1, \dots, \tau_{\ell_i})$	(Sum Type)
$\&p$	(Reference pattern)	$\&\tau$	(Reference Type)
$\_$	(Wildcard)	$\Gamma ::= \{x_1 : \tau_1; \dots; x_n : \tau_n\}$	(Type Environment)
$x$	(Variable)	Values	
$(p_1 \mid p_2)$	(Disjunction)	$v ::= c \in \mathcal{C}$	(Constant)
		$\langle v_1, \dots, v_n \rangle$	(Product)
		$K(v_1, \dots, v_n)$	(Constructor)
		$\&v$	(Reference)
		$\sigma ::= \{x_1 \mapsto v_1; \dots; x_n \mapsto v_n\}$	(Value Env.)
Matching			
$m ::= \text{Match} \{p_1 \mid \dots \mid p_n\}$	(Matching)		

Figure 5: Our language of pattern and their types

ANY $\_ \triangleright v \rightarrow \emptyset$	VAR $x \triangleright v \rightarrow \{x \mapsto v\}$	CONSTANT $c \triangleright c \rightarrow \emptyset$	REFERENCE $\frac{p \triangleright v \rightarrow \sigma}{\&p \triangleright \&v \rightarrow \sigma}$
TUPLE $\frac{\forall i, p_i \triangleright v_i \rightarrow \sigma_i \quad \forall j \neq i, \sigma_i \cap \sigma_j = \emptyset}{\langle \bar{p}_i \rangle \triangleright \langle \bar{v}_i \rangle \rightarrow \bigcup \sigma_i}$	CONSTRUCTOR $\frac{\forall i, p_i \triangleright v_i, \sigma_i \quad \forall j \neq i, \sigma_i \cap \sigma_j = \emptyset}{K(\bar{p}_i) \triangleright K(\bar{v}_i) \rightarrow \bigcup \sigma_i}$	ALTL $\frac{p_1 \triangleright v \rightarrow \sigma}{(p_1 \mid p_2) \triangleright v \rightarrow \sigma}$	
ALTR $\frac{p_1 \not\triangleright v \quad p_2 \triangleright v \rightarrow \sigma}{(p_1 \mid p_2) \triangleright v \rightarrow \sigma}$	MATCHING $\frac{p_i \triangleright v \rightarrow \sigma \quad \forall j < i, p_j \not\triangleright v}{\text{Match} \{p_1 \mid \dots \mid p_n\} \triangleright v \rightarrow i, \sigma}$		

Figure 6: Semantic of patterns –  $p \triangleright v \rightarrow \sigma$ 

- We enforce that there is no shadowing: variables must be bound only once, as asserted by the side-condition in the TUPLE and CONSTRUCTOR rules.
- Alternatives are left-leaning: we first try to match with the left branch (ALTL rule) before trying the right branch (ALTR rule).

We also define the matching judgment in rule MATCHING:  $\text{Match} \{p_1 \mid \dots \mid p_n\} \triangleright v \rightarrow i, \sigma$  which behaves as the normal judgment, but additionally returns the index of the branch which was matched. In a full language, it would then trigger the evaluation of the body of the branch in question.

**Example 4** (Interval example, cont'). For the following matching problem in which we provide pattern names:

$$\text{Match} \left\{ \begin{array}{l} | (p_0) \text{ None} \\ | (p_1) \text{ Some}(\&(x_1, x_2)) \end{array} \right\}$$

- The *pattern*  $p_0$  will match the value **None**, which will be denoted by  $p_0 \triangleright \text{None} \rightarrow 0, \emptyset$ ;
- The second pattern,  $p_1$ , will match any **Interval** inhabitant of the form  $\text{Some}(\&(v, v'))$ , with  $v, v'$  two integers. For instance  $p_1 \triangleright \text{Some}(\&(1756, 1791)) \rightarrow 1, \{x_1 \mapsto 1756; x_2 \mapsto 1791\}$ .

### 3.2 Typing

Fig. 7 defines the typing judgment  $\vdash p : \tau \rightarrow \Gamma$  which means that pattern  $p$  has type  $\tau$  and binds variables present in environment  $\Gamma$ . Typing follows the semantics closely:

$$\begin{array}{c}
\text{ANY} \\
\frac{}{\vdash \_ : \tau \rightarrow \emptyset} \\
\\
\text{VAR} \\
\frac{}{\vdash x : \tau \rightarrow \{x : \tau\}} \\
\\
\text{CONSTANT} \\
\frac{c \in T}{\vdash c : T \rightarrow \emptyset} \\
\\
\text{REFERENCE} \\
\frac{\vdash p : \tau \rightarrow \Gamma}{\vdash \&p : \&\tau \rightarrow \Gamma} \\
\\
\text{TUPLE} \\
\frac{\forall i; \vdash p_i : \tau_i \rightarrow \Gamma_i \quad \forall i, j; \Gamma_i \cap \Gamma_j = \emptyset}{\vdash \langle \overline{p_i} \rangle : \prod \tau_i \rightarrow \bigcup \Gamma_i} \\
\\
\text{ALT} \\
\frac{\forall i; \vdash p_i : \tau \rightarrow \Gamma}{\vdash p_1 \mid p_2 : \tau \rightarrow \Gamma} \\
\\
\text{CONSTRUCTOR} \\
\frac{\exists i_0, K = K_{i_0} \quad \forall j; \vdash p_j : \tau_{i,j} \rightarrow \Gamma_j \quad \forall i, j; \Gamma_i \cap \Gamma_j = \emptyset}{\vdash K(\overline{p_j}) : \sum K_i(\overline{\tau_{i,j}}) \rightarrow \bigcup \Gamma_j} \\
\\
\text{MATCHING} \\
\frac{\forall i; \vdash p_i : \tau \rightarrow \Gamma}{\vdash \text{Match} \{p_1 \mid \dots \mid p_n\} : \tau \rightarrow \Gamma}
\end{array}$$

Figure 7: Typing of patterns –  $\vdash p : \tau \rightarrow \Gamma$ 

- As before, the bound typing environment  $\Gamma$  is populated by variables through the VAR rule.
- We enforce that there is no shadowing in the TUPLE and CONSTRUCTOR rules.
- Bound environments must be identical in all branches, as enforced in the ALT rule.

**Example 5** (Pattern types for Interval example).

$$\vdash p_1 : \text{Interval} \rightarrow \emptyset \qquad \vdash p_2 : \text{Interval} \rightarrow \{x_1 : u64; x_2 : u64\}$$

## 4 Target Language and Memory Representations

We now describe the target environment of our compilation procedure. In the context of this article, part of the target environment is considered universal: what kind of structures can be represented in memory (such as machine words, pointers, ...) and how to manipulate and compute them. Another part of the target is specific and defines a so-called “memory representation”: a choice of how language values are represented in memory. We formally state the notion of *valid* memory representations which is crucial for our compilation procedure, along with some examples.

### 4.1 Memory values

In order to reason about values in memory, we consider an abstraction of memory contents, whose grammar is presented in Fig. 8a. Memory values can be arbitrarily long machine words, denoted by  $w$ ;  $n$ -bit wide pointers to another memory value, denoted by  $\&_n r$ ; or contiguous arrays of memory values with alignment  $a$ , denoted by  $[r_0, \dots, r_n]_a$ .

**Example 6** (Memory value of *Some*(5) in OCaml). In OCaml, *Some*(5) of type `int option` (an optional integer) is represented in memory by a pointer to a block (a contiguous array) containing a tag followed by the integer 5. The tag in this case is 0. The integer 5 is represented as  $\overline{1011}^2$  due to the pointer flag Minsky and Madhavapeddy (2021). The resulting representation of *Some*(5) can thus be written as  $\&_{64} [0, \overline{1011}^2]_{64}$ .

### 4.2 Computing on memory values

#### 4.2.1 Decision trees

From a matching problem, we will produce our target under the form of *decision trees*, defined in Fig. 8c. Decision trees are composed of nodes of the form `switch(e) { C }`, consisting of an expression  $e$  computing the memory value under scrutiny – also called the *switch discriminant* – and of a list of cases  $\mathcal{C}$ . Leaves of a decision tree can either be `success(j, Θ)`, which successfully returns a branch index  $j$  and a set of binding

<p><b>Memory values</b></p> $r ::= w \quad (\text{Word})$ $  \&_n r \quad (n\text{-bit wide reference})$ $  [\bar{r}_i]_a \quad (a\text{-aligned contiguous array})$ <p>(a) Memory values</p> <p><b>Tags</b></p> $t ::= \top \quad (\text{Default tag})$ $  c \quad (\text{Constant tag})$ $  K \quad (\text{Constructor tag})$ <p><b>Contexts</b></p> $h ::= \square \quad (\text{Hole})$ $  \langle \_, \dots, \_, h, \_, \dots, \_ \rangle \quad (\text{Product context})$ $  K(\_, \dots, \_, h, \_, \dots, \_) \quad (\text{Constr. context})$ $  \&h \quad (\text{Reference context})$ <p>(b) Tags and Contexts</p>	<p><b>Target expressions</b></p> $e ::= \Delta \quad (\text{Main Discriminant})$ $  * e \quad (\text{Dereferencing})$ $  e.\text{offset} \quad (\text{Array access})$ $  \dots \quad (\text{Other operations})$ <p><b>Binding expressions</b></p> $\Theta ::= \{x_i \mapsto e_i\}$ <p><b>Switch cases</b></p> $\mathcal{C} ::= w \mapsto \mathcal{T} \cdot \mathcal{C} \quad (\text{Regular case})$ $  \top \mapsto \mathcal{T} \quad (\text{Default case})$ $  \emptyset \quad (\text{No Default case})$ <p><b>Decision trees</b></p> $\mathcal{T} ::= \text{switch}(e) \{ \mathcal{C} \} \quad (\text{Decision node})$ $  \text{success}(j, \Theta) \quad (\text{Branch } j \text{ with bindings } \Theta)$ $  \text{unreachable} \quad (\text{Unreachable leaf})$ <p>(c) Decision trees</p>
---	--

Figure 8: The target languages

$\frac{\text{SUCCESS} \quad \forall (x \mapsto e) \in \Theta, \quad e[r] \hookrightarrow r_x}{\text{success}(j, \Theta)[r] \hookrightarrow j, \{x \mapsto r_x\}_{x \in S}}$	$\frac{\text{DEFAULT} \quad \mathcal{T}[r] \hookrightarrow j, \sigma}{\text{switch}(e) \{ \top \mapsto \mathcal{T} \}[r] \hookrightarrow j, \sigma}$	$\frac{\text{CASET} \quad \mathcal{T}[r] \hookrightarrow j, \sigma \quad e[r] \hookrightarrow w}{\text{switch}(e) \{ w \mapsto \mathcal{T} \cdot l \}[r] \hookrightarrow j, \sigma}$
$\frac{\text{CASEF} \quad \text{switch}(e) \{ l \}[r] \hookrightarrow j, \sigma \quad e[r] \hookrightarrow w' \neq w}{\text{switch}(e) \{ w \mapsto \mathcal{T} \cdot l \}[r] \hookrightarrow j, \sigma}$		

Figure 9: Decision tree evaluation

expressions  $\Theta$ , or `unreachable`.  $\Theta$  can be understood as a set of let-binders which are necessary for the later evaluation of the successful branch. At toplevel, decision trees take as input the *main discriminant* of the whole matching, designated by  $\Delta$ . The discriminant of a (sub)-switch is an expression manipulating memory values mentioning this main discriminant. In addition, expressions are target-dependent and can include dereferencing denoted by  $*e$ , word accesses denoted by  $e.\text{offset}$  (where the offset is statically known), arithmetic and logical operations, etc.

A decision tree  $\mathcal{T}$  applied to an input memory value  $r$  reduces to a result  $j, \sigma$  where  $j$  is a branch index and  $\sigma$  is a binding environment. This judgment is written  $\mathcal{T}[r] \hookrightarrow j, \sigma$  and is defined in Fig. 9. We assume that a similar judgment reduces expressions to memory values, and is denoted by  $e[r] \hookrightarrow r'$ . The rule `SUCCESS` evaluates a successful leaf by evaluating each expression contained in the binders  $\Theta$ . The three other rules `DEFAULT`, `CASET` and `CASEF` evaluate a switch node by finding the first case matching the discriminant: either the first branch whose left-hand side is equal to the discriminant, or the wildcard (if present). There are no rules for evaluating an empty branch or an “`unreachable`” leaf: as mentioned previously, we consider all source-level language patterns to be exhaustive. We will prove that such cases are indeed unreachable for all emitted trees.

**Example 7** (Decision tree for `Interval`). For illustration purposes, we translate the decision tree shown

in Fig. 2, for the Rust version:

$$\text{switch}(\Delta) \left\{ \begin{array}{l} 0 \rightarrow \text{success}(0, \emptyset) \\ 1 \rightarrow \text{success}(1, \sigma) \end{array} \right\}$$

In this decision tree we identify input **Intervals** of the form  $\text{Some}(\&(x_1, x_2))$ . At toplevel, the decision tree checks whether the considered memory value  $\Delta$  is 0 or 1, which is enough to decide whether the Interval is **None** or not. The *binding environment*  $\sigma$  for “success at branch 1” maps, for instance,  $x_1$  to the left value of the interval.

#### 4.2.2 Tags and Contexts

To fully describe memory representations and later our compilation scheme, we need two additional tools: tags and contexts. Their grammar is given in Fig. 8b.

*Tags* are the simplest description of the possible variants of a type. They can be either a constant of a primitive type (integer, float, char...), a constructor of a sum type or a wildcard  $\top$ . Other types such as references and product types do not have associated tags, as they can’t give rise to alternatives in patterns by themselves.

*Contexts* are used to deconstruct values and types. Contexts are a subclass of patterns where all sub-patterns are wildcards, except one, which is a hole denoted by  $\square$ . We will often use contexts to precisely describe the structure around subterms bound by a variable in a pattern.

*Focusing* is a set of operations related to contexts.  $\mathbf{focus}(h, v)$  returns the sub-value of  $v$  at the position of  $\square$  in  $h$ .  $\mathbf{focus}(h, \tau)$  returns the type of the subterm at the position of  $\square$  for values of type  $\tau$ . Both judgments are fully defined in Appendix B. For now, we only give the intuition through Example 8.

**Example 8** (Contexts for **Intervals**). Let us consider the pattern  $p = \text{Some}(\&(x_1, x_2))$ . The context corresponding to  $x_2$  in  $p$  is  $h = \text{Some}(\&(\_, \square))$ . We thus have:

$$\begin{aligned} \mathbf{focus}(h, \text{Interval}) &= u64 \\ \mathbf{focus}(h, \text{Some}(\&(u1756, u1791))) &= 1791 \end{aligned}$$

### 4.3 Memory representation

The compilation process depicted in Section 5 will be *parametrized* by a memory representation. The objective of this section is to define the required ingredients for this parametrization, and also define what would be a *valid* representation.

A *memory representation*  $\mathfrak{R}$  is a triple  $\mathfrak{R} = (\text{REPR}^\bullet, \text{FROG}^\bullet, \text{KNIT}^\bullet)$ . Indeed, giving a function  $\text{REPR}^\bullet : \text{VAL} \rightarrow \text{MEMVAL}$  computing memory values from values (Section 4.1) is not sufficient: decision trees should be produced with their target-dependent expressions inside switches. For instance, in Example 7, the second switch should come with a representation-dependent way to compute  $(*\Delta).2$ . Moreover, we also need a way to generate alternatives and their associated selection code. We propose two additional methods for which an example is given in Section 4.4. For ease of reading, operations specific to a given representation are denoted with a full dot.

- $\text{FROG}_{\tau_\Delta}^\bullet : \mathcal{H} \rightarrow (\overline{t \mapsto \mathcal{T}}) \rightarrow \mathcal{T}$  creates representation-specific code to choose between several alternatives in a given type<sup>4</sup>. It takes a context  $h$  indicating a subterm of the current pattern (the subterm of the main discriminant  $\Delta$  under scrutiny) and a list of tags and their associated decision trees  $(\overline{t_k, \mathcal{T}_k})$ . It returns a

$$\text{complete decision tree } \mathcal{T} = \text{switch}(e) \left\{ \begin{array}{l} w_1 \mapsto \mathcal{T}_1 \\ \vdots \\ w_n \mapsto \mathcal{T}_n \end{array} \right. \quad \text{where } e \text{ is a representation-dependent expression that}$$

characterizes the subterm and each  $w_k$  is one of its possible runtime values. The semantics of this tree is “expression  $e$  computes a value that exactly characterizes the object inside the hole  $h$ ; we choose a subtree through which to continue according to this value”.

Each tag  $t_k$  is either a possible variant for the destructured type  $\mathbf{focus}(h, \tau_\Delta)$  that encompasses all values of this variant (see Section 4.2.2) or the special default tag  $\top$  (if the image of the  $\mathbf{focus}(h, \tau)$  function does not cover all possible values of the type  $\tau$ ).

<sup>4</sup>For crochet enthusiasts, *Frogging* is the action of undoing the stitches: rip it, ripit, ribbit, ribbit, ...

- $\text{KNIT}_{\tau_{\Delta}}^{\bullet} : \mathcal{H} \rightarrow e$ .  $\text{KNIT}^{\bullet}$  creates representation-specific code to build the memory value representing a subterm of the main discriminant. It takes a context that indicates the subterm under scrutiny and returns an expression computing the memory value representing this subterm. As shown in Example 7, these memory values might not be directly accessible from the parent memory value and might need to be reconstructed based on various pieces present in the parent memory value, hence the need for such a function.

Naturally, there are constraints over what determines a reasonable choice of  $\text{REPR}^{\bullet}$ ,  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$ . The main idea is that  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$  should be compatible with  $\text{REPR}^{\bullet}$ : the expression produced by  $\text{KNIT}^{\bullet}$  should evaluate to a final memory value which is the same as the one built by  $\text{REPR}^{\bullet}$ , and the decision tree produced by  $\text{FROG}^{\bullet}$  should properly distinguish between different variants based on the memory values produced by  $\text{REPR}^{\bullet}$ .

**Definition 1** (Valid Representation). *We say that a triple  $\mathfrak{R} = (\text{REPR}^{\bullet}, \text{FROG}^{\bullet}, \text{KNIT}^{\bullet})$  is a valid representation if, and only if, the following holds.*

*Let  $(t_1, \mathcal{T}_1) \dots (t_n, \mathcal{T}_n)$  be the branches under consideration,  $h$  a context,  $\tau$  the type associated with the tags and  $v$  a value of type  $\tau$ .*

*We define the auxiliary function  $\text{Tag2Pat} : \text{Tags} \mapsto \text{PAT}$  s.t.*

$$\text{Tag2Pat}(c) = c \qquad \text{Tag2Pat}(K) = K(-, \dots, -) \qquad \text{Tag2Pat}(\top) = -$$

*Then we should have:*

$$\text{KNIT}_{\tau}^{\bullet}(h)[\text{REPR}_{\tau}^{\bullet}(v)] \leftrightarrow \text{REPR}_{\text{focus}(h, \tau)}^{\bullet}(\text{focus}(h, v))$$

*and:*

$$\begin{aligned} & \text{FROG}_{\tau}^{\bullet}\left(h, \overline{(t_k, \mathcal{T}_k)}\right)[\text{REPR}_{\tau}^{\bullet}(v)] \leftrightarrow j, \sigma \\ & \exists k. \text{Tag2Pat}(t_1) \mid \dots \mid \text{Tag2Pat}(t_n) \updownarrow v \rightarrow k \wedge \mathcal{T}_k[\text{REPR}_{\tau}^{\bullet}(v)] \leftrightarrow j, \sigma \end{aligned}$$

## 4.4 A first example: the boxed representation

### 4.4.1 Memory values

The boxed representation maps values of all types to uniform, word-wide memory values. The complete definition is given in Fig. 10.

Constant and reference values can always fit into a single memory word and are thus stored *unboxed*. (Numeric) constants are stored as words in a type-appropriate representation/encoding. The memory value of a reference is a pointer to the memory value of the referenced value. Access is done by dereferencing the pointer. Tuple and constructor values may not fit into a single word, hence the need for *blocks*, which are word-aligned contiguous arrays. More precisely, tuples are represented by pointers to blocks in which the  $k$ -th word is a memory value representing the  $k$ -th field of the tuple. The representation of constructor values is similar, with the index of the constructor among all possible constructors for the considered type stored in the first word of the block and the  $k$ -th field in the  $k + 1$ -th word. Accessing the  $k$ -th field of such a value consists in dereferencing the pointer and accessing the  $k$ -th or  $k + 1$ -th word of the block.

**Example 9** (Boxed representation, memory values). Let us consider the type  $\tau_0 = \text{None} + \text{Some}(A + B + C(\text{Int}_{32}))$ , which represent an optional value, itself containing three possible cases. Here are a few examples of values according to the boxed representation:

$$\begin{aligned} \text{REPR}_{\tau_0}^{\bullet}(\text{None}) &= \&_{64}[0]_{64} & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(A)) &= \&_{64}[1, \&_{64}[0]_{64}]_{64} \\ \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(B)) &= \&_{64}[1, \&_{64}[1]_{64}]_{64} & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(C(n))) &= \&_{64}[1, \&_{64}[2, n]_{64}]_{64} \end{aligned}$$

$\begin{aligned} \text{REPR}_{\top}^{\bullet} &: c \mapsto c \\ \text{REPR}_{\&\tau}^{\bullet} &: \&v \mapsto \&_{64}(\text{REPR}_{\tau}^{\bullet}(v)) \\ \text{REPR}_{\prod \tau_i}^{\bullet} &: \langle \bar{v}_i \rangle \mapsto \&_{64} \left[ \overline{\text{REPR}_{\tau_i}^{\bullet}(v_i)} \right]_{64} \\ \text{REPR}_{\sum K_i(\bar{\tau}_{i,j})}^{\bullet} &: K_i(\bar{v}_j) \mapsto \&_{64} \left[ i, \overline{\text{REPR}_{\tau_{i,j}}^{\bullet}(v_j)} \right]_{64} \end{aligned}$	$\begin{aligned} f(\square, e) &= e \\ f(\&h, e) &= f(h, *e) \\ f(\langle h_k \rangle, e) &= f(h, *e.k) \\ f(K(h_k), e) &= f(h, *e.k + 1) \\ \text{KNIT}^{\bullet}(h) &= f(h, \Delta) \end{aligned}$
(a) Memory values	(b) Knitting

**Frogging, constant type – focus**  $(h, \tau) = T$

$$\text{FROG}_{\tau}^{\bullet}(h) \left( \overline{(c_k, \mathcal{T}_k)} \cdot (\top, \mathcal{T}_*) \right) = \text{switch}(\text{KNIT}^{\bullet}(h)) \begin{cases} \overline{c_k} \mapsto \overline{\mathcal{T}_k} \\ \top \mapsto \mathcal{T}_* \end{cases}$$

**Frogging, sum type – focus**  $(h, \tau) = \sum K_i(\bar{\tau})$

$$\begin{aligned} \text{FROG}_{\tau}^{\bullet}(h) \left( \overline{(K_k, \mathcal{T}_k)} \right) &= \text{switch}((\text{*KNIT}^{\bullet}(h)).0) \begin{cases} \overline{k} \mapsto \overline{\mathcal{T}_k} \end{cases} \quad \text{with } \{\overline{K_k}\} \text{ a complete signature} \\ \text{FROG}_{\tau}^{\bullet}(h) \left( \overline{(K_{i_k}, \mathcal{T}_k)} \cdot (\top, \mathcal{T}_*) \right) &= \text{switch}((\text{*KNIT}^{\bullet}(h)).0) \begin{cases} \overline{i_k} \mapsto \overline{\mathcal{T}_k} \\ \top \mapsto \mathcal{T}_* \end{cases} \end{aligned}$$

(c) Frogging

Figure 10: Boxed representation with 64-bit words

#### 4.4.2 Representation (Knit&Frog)

To distinguish between constructors,  $\text{FROG}^{\bullet}$  emits code which dereferences the pointer and accesses the first word of the block, and  $\text{switch}()$  on it.

The boxed representation has the remarkable property of being *composable*: memory value representing a value always contains the memory values of its subterms. Concretely, almost all functions manipulating the representation can be defined by induction without much inspection of the sub-types. For instance,  $\text{REPR}^{\bullet}$ , defined in Fig. 10a, is a direct induction over the structure of values. Similarly, knitting is defined through an auxiliary function  $f$  shown in Fig. 10b which takes a context, and an expression computing the parent value of that context. It then inductively accumulates the expression by deconstructing the context.  $\text{KNIT}^{\bullet}$  is a restriction of  $f$  to the main discriminant  $\Delta$ . Furthermore, accesses to constructors of values are uniform across types and variants.  $\text{FROG}^{\bullet}$ , defined in Fig. 10c, directly uses  $\text{KNIT}^{\bullet}$  to reach the constructor of a memory value and always yields a single switch to distinguish between all variants of a given type.

**Example 10** (Boxed representation, cont'). Examples of  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$  on values of type  $\tau_0$  of Example 9 are given in Fig. 11. For  $\text{KNIT}^{\bullet}$ , to compute the subterm induced by a context, we simply access the piece of memory contained at its position with simple dereferencing and indexed access. The decision tree  $\text{FROG}_{\tau_0}^{\bullet}(\text{Some}(\square))((B, \mathcal{T}_1) \cdot (\top, \mathcal{T}_*))$  should distinguish *whether the hole* is a  $B$  or anything else. It is thus sufficient to look at the discriminant  $(\text{*}\Delta).1.0$  (Example 9). In the general case  $\text{FROG}^{\bullet}$  accesses the relevant tag and distinguish cases in a uniform manner.

The boxed memory values and its  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$  operations closely mimic pattern matching on values in the surface language, we thus state the following result:

**Lemma 1.** The boxed representation is a valid representation.

As we have seen, thanks to its very uniform definition the boxed representation is easy to state and manipulate.  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$  can be defined in term of each other, and a call to  $\text{FROG}^{\bullet}$  always yields a single switch. This is rarely the case in more complex representations. Naturally, this simplicity induces

**Knitting**

$$\text{KNIT}_{\tau_0}^{\bullet}(\text{Some}(\square)) = (*\Delta).1 \qquad \text{KNIT}_{\tau_0}^{\bullet}(\text{Some}(C(\square))) = *((*\Delta).1).1$$

**Frogging**

$$\begin{aligned} \text{FROG}_{\tau_0}^{\bullet}(\square) ((None, \mathcal{T}_0) \cdot (\text{Some}, \mathcal{T}_1)) &= \text{switch} ((*\Delta).0) \begin{cases} 0 \mapsto \mathcal{T}_0 \\ 1 \mapsto \mathcal{T}_1 \end{cases} \\ \text{FROG}_{\tau_0}^{\bullet}(\text{Some}(\square)) ((B, \mathcal{T}_1) \cdot (\top, \mathcal{T}_*)) &= \text{switch} ((*(*\Delta).1).0) \begin{cases} 1 \mapsto \mathcal{T}_1 \\ \top \mapsto \mathcal{T}_* \end{cases} \\ \text{FROG}_{\tau_0}^{\bullet}(\text{Some}(C(\square))) ((42, \mathcal{T}_{42}) \cdot (\top, \mathcal{T}_*)) &= \text{switch} ((*(*\Delta).1).1) \begin{cases} 42 \mapsto \mathcal{T}_{42} \\ \top \mapsto \mathcal{T}_* \end{cases} \end{aligned}$$

Figure 11: Memory values, knitting and frogging for values of type  $\tau_0$ 

numerous inefficiencies: Every constructor is represented as a pointer to a block, including argument-less variants such as *None* which could easily be represented as unboxed integers. Frogging a value with the context  $\text{Some}(\square)$  thus requires dereferencing two pointer, while only one might be necessary for a more packed representation. Mainstream languages with ADTs such as OCaml, Haskell, or Rust apply such optimization to their representation with various means of differentiating between constructors and values. Section 6 showcases more optimized representation in order to demonstrate our compilation process in a more realistic setting.

## 5 Representation-dependent compilation

We now define our compilation scheme. In the previous sections, we have defined pieces of the target decision trees we want to emit. In isolation, those pieces are easy to deduce from simple patterns. The difficulty, however, comes with composition: how to handle nested patterns and how to expose back the values bound through variables in a way that is compatible with the chosen memory representation. This composition is precisely the object of our compilation procedure.

Our procedure is inspired by Maranget (2008). The novelty here is that this procedure is now parameterized by a *valid representation* and creates decision trees manipulating memory values directly. Our contribution is the careful split between the representation-dependent elements, which are parameterized over and delegated to the  $\text{KNIT}^{\bullet}$  and  $\text{FROG}^{\bullet}$  functions as described in the previous sections, and the representation-independent aspects which are described here. Additionally, we add handling of the bound environment, necessary for the rest of the compilation process.

Let us start by stating our goal: we aim to compile a matching  $m$  to a *decision tree*  $\mathcal{T}$  that takes a memory value  $r$  as input and evaluates to an output (denoted by  $\mathcal{T}[\overline{r}_i] \hookrightarrow j, \sigma$ , see Section 4.2.1) consistent with the matching judgment on the source values.

More formally, given a (source-level) value  $v$  of type  $\tau$ , we want :

$$m \triangleright v \rightarrow j, \sigma \iff \mathcal{T}[\text{REPR}_{\tau}^{\bullet}(v)] \hookrightarrow j, \{x \mapsto \text{REPR}_{\tau_x}^{\bullet}(y) \mid (x \mapsto y) \in \sigma\}$$

In the rest of this section, we consider as an implicit ambient parameter a *valid representation*  $\mathfrak{R} = (\text{REPR}^{\bullet}, \text{FROG}^{\bullet}, \text{KNIT}^{\bullet})$ . As stated before, operations specific to a given representation are denoted with a full dot:  $\text{REPR}^{\bullet}$ . Operations which are *parameterized* by the representation are denoted with an empty dot:  $\text{COMPILE}_{\tau_{\Delta}}^{\circ}$ . Operations without any dot are representation-agnostic.



## 5.1 Algorithm

$\text{COMPILE}_{\tau_{\Delta}}^{\circ}(\mathcal{P})$  compiles the given patterns to a decision tree. It takes two parameters: the type of main discriminant  $\tau_{\Delta}$ , and the *pattern matrix*  $\mathcal{P}$ , which gathers all the necessary ingredients. In its most general form, the pattern matrix is composed as so:

$$\mathcal{P} = \left( \begin{array}{cccc|c} h_1 & \dots & h_i & \dots & h_n & j^1, s^1 \\ p_1^1 & & p_i^1 & & p_n^1 & \vdots \\ \vdots & & \vdots & & \vdots & \vdots \\ p_1^\ell & \dots & p_i^\ell & \dots & p_n^\ell & j^\ell, s^\ell \\ \vdots & & \vdots & & \vdots & \vdots \\ p_1^m & & p_i^m & & p_n^m & j^m, s^m \end{array} \right)$$

where:

- The input matching problem is encoded in rows: row  $k$  represents a case of the match ( $(p_1^k, \dots, p_n^k)$ ). These rows end with a main branch index ( $j^k$ ) and a binding environment ( $s^k$ ). Importantly, we consider a *vector* of discriminants instead of a single one as done previously, allowing us to split product patterns and choose which one to inspect first.
- Each column represents the different parts of the discriminant under consideration. Their headers  $h_i$  are contexts that indicate which parts of the main discriminant this column is matching against, as a subterm of the main discriminant  $\Delta$ . For instance, on inspecting an optional integer, we might discover that it is a *Some(...)* and want to explore the subterm in context *Some(□)* of the discriminant  $\Delta$ .

$\text{COMPILE}^{\circ}$  proceeds by recursively emitting Switch nodes, and reducing the pattern matrix until exhaustion. It has five cases, depicted in Fig. 12.

- If  $\mathcal{P}$  is empty, the branch is *unreachable* as there are no patterns to inspect the discriminant.
- If the first row consists only of wildcards, the matching can only succeed.
- If there exists a variable pattern  $p_i^\ell$  and the previous rows contain no variable patterns, we introduce  $p_i^\ell$  in the binding environment.
- If there exists an “or”-pattern  $p_i^\ell = p_1 \mid p_2$  and the previous rows contain no “or”-patterns, split it.
- Otherwise, we must build an actual switch node. This case is more involved and detailed in the rest of this section.

The aim of the switch node is to inspect the head constructor of the value contained in one of the discriminants against the patterns of its column in  $\mathcal{P}$ , and to branch to the sub-trees that perform the remaining computations (on other columns and nested values). Note that we assume  $\mathcal{P}$  does not contain any variable or “or”-patterns, as those are removed by the two previous cases. The construction of a switch node (last case of Fig. 12) is done in four steps, detailed as follows.

### 1. Pick a column to inspect

We first pick an arbitrary column  $i$  such that  $p_i^1 \neq \_$  to generate a node switching on that column. This pick is done by the function `PICKCOLUMN`, which takes the pattern matrix  $\mathcal{P}$ , and returns a column index. The choice of column can be done through heuristics as described by Maranget (2008); Scott and Ramsey (2000). It does not affect our safety proof.

### 2. Identify all the necessary branches for the chosen column

Now that we have chosen a column, we need to identify the different branches of our switch node. Let us consider the  $i$ -th column: it is associated with a discriminant designated by a context  $h_i$ . We can obtain its type thanks to focusing:  $\tau_i = \mathbf{focus}(h_i, \tau_{\Delta})$ . The set of necessary branches depends directly on the type of the column  $\tau_i$ : if the discriminant represents a reference or tuple value, then only one possible branch is needed. Otherwise, if  $\tau_i$  is a primitive or sum type, several branches are needed. More precisely, the generated switch node should include a branch for each tag that appears as the head constructor of a pattern in column  $i$ , with an additional default branch if those constructors do not encompass all possible values of this type. The auxiliary representation-independent function `GETTAGS $_{\tau_i}$`  defined below

$$\begin{aligned}
 & \text{COMPILE}_{\tau_{\Delta}}^{\circ}(\emptyset) = \text{unreachable} && \text{(No pattern case)} \\
 & \text{COMPILE}_{\tau_{\Delta}}^{\circ} \left( \begin{array}{c|c} h_1 & \dots & h_n \\ \hline - & \dots & - \\ \vdots & \ddots & \vdots \\ p_1^m & \dots & p_n^m \end{array} \middle| \begin{array}{c} j^1, s^1 \\ \vdots \\ j^m, s^m \end{array} \right) = \text{success}(j^1, s^1) && \text{(Wildcard case)} \\
 & \text{COMPILE}_{\tau_{\Delta}}^{\circ} \left( \begin{array}{c|c} h_1 & h_i & h_n \\ \hline p_1^1 & p_i^1 & p_n^1 \\ \vdots & \vdots & \vdots \\ p_1^\ell & \dots & p_n^\ell \\ \vdots & \vdots & \vdots \\ p_1^m & p_i^m & p_n^m \end{array} \middle| \begin{array}{c} j^1, s^1 \\ \vdots \\ j^\ell, s^\ell \\ \vdots \\ j^m, s^m \end{array} \right) = \text{COMPILE}_{\tau_{\Delta}}^{\circ} \left( \begin{array}{c|c} h_1 & h_i & h_n \\ \hline p_1^1 & \dots & p_n^1 \\ \vdots & \vdots & \vdots \\ p_1^\ell & \dots & p_n^\ell \\ \vdots & \vdots & \vdots \\ p_1^m & \dots & p_n^m \end{array} \middle| \begin{array}{c} j^1, s^1 \\ \vdots \\ j^\ell, s^\ell \cup s_x \\ \vdots \\ j^m, s^m \end{array} \right) \\
 & \hspace{15em} \text{where } s_x = \{x \rightarrow \text{KNIT}_{\tau_{\Delta}}^{\bullet}(h_i)\} && \text{(Variable case)} \\
 & \text{COMPILE}_{\tau_{\Delta}}^{\circ} \left( \begin{array}{c|c} h_1 & h_i & h_n \\ \hline p_1^1 & \dots & p_n^1 \\ \vdots & \vdots & \vdots \\ p_1^\ell & \dots & p_n^\ell \\ \vdots & \vdots & \vdots \\ p_1^m & \dots & p_n^m \end{array} \middle| \begin{array}{c} j^1, s^1 \\ \vdots \\ j^\ell, s^\ell \\ \vdots \\ j^m, s^m \end{array} \right) = \text{COMPILE}_{\tau_{\Delta}}^{\circ} \left( \begin{array}{c|c} h_1 & h_i & h_n \\ \hline p_1^1 & \dots & p_n^1 \\ \vdots & \vdots & \vdots \\ p_1^\ell & \dots & p_n^\ell \\ \vdots & \vdots & \vdots \\ p_1^m & \dots & p_n^m \end{array} \middle| \begin{array}{c} j^1, s^1 \\ \vdots \\ j^\ell, s^\ell \\ \vdots \\ j^m, s^m \end{array} \right) \\
 & \hspace{15em} \text{(Or case)} \\
 & \text{COMPILE}_{\tau_{\Delta}}^{\circ}(\mathcal{P}) = \begin{cases} i \leftarrow \text{PICKCOLUMN}(\mathcal{P}) & (1) \\ \text{Let } h_i \text{ the } i\text{-th context in } \mathcal{P} \text{ and } \tau_i = \mathbf{focus}(h_i, \tau_{\Delta}) & (2) \\ \text{Tags} \leftarrow \text{GETTAGS}_{\tau_i}(i, \mathcal{P}) & (3) \\ \mathcal{C} \leftarrow \bigsqcup_{tag \in \text{Tags}} (t, \text{COMPILEBRANCH}_{\tau_{\Delta}}^{\circ}(\mathcal{P})(i, tag)) & (3) \\ \text{FROG}_{\tau_{\Delta}}^{\bullet}(h_i)(\mathcal{C}) & (4) \end{cases} && \text{(Switch case)}
 \end{aligned}$$

 Figure 12: The compilation procedure –  $\text{COMPILE}_{\tau_{\Delta}}^{\circ}(\mathcal{P})$ 

enumerates the *tags* (as defined in Section 4.2.2) associated with the necessary sub-trees. It works by collecting all the tags present in the pattern matrix at the column  $i$ . For constant types, we always add a “catch-all” case  $\top$ . For sum types, we add it only if the collected tags do not cover all the constructors in the type  $\tau_i$ . Finally, for all other cases, as described before, there is only one case corresponding to the generic  $\top$  tag.

$$\text{GETTAGS}_T(i, \mathcal{P}) = \{c \mid \exists j. p_i^j = c\} ++ \top \quad \text{where } ++ \text{ is concatenation}$$

$$\text{GETTAGS}_{\Sigma}(i, \mathcal{P}) = \{K \mid \exists j. p_i^j = K(\dots)\} ++ \begin{cases} \{\top\} & \text{if } \exists K'. \forall j. p_i^j = K'(\dots) \wedge K' \neq K \\ \{\} & \text{otherwise} \end{cases}$$

$$\text{GETTAGS}_{\_}(i, \mathcal{P}) = \top$$

### 3. Compile the decision sub-trees corresponding to each branch

Now that we have identified the tag corresponding to each branch, we can compute the corresponding subtree for each tag. Each subtree is the result of the compilation of a new pattern matrix, as done by

Inputs		Outputs	
$\tau$	$tag$	$\mathcal{H}$	$f$
$T$	$c$	$()$	$c \mapsto ()$ $c' \mapsto \emptyset \quad c' \neq c$ $- \mapsto ()$
$T$	$\top$	$()$	$c \mapsto \emptyset$ $- \mapsto ()$
$\&\tau$	$\top$	$(h[\&\square])$	$\&q \mapsto (q)$ $- \mapsto (-)$
$\prod_{1 \leq k \leq \ell} \tau_k$	$\top$	$(h[h_0] \ \dots \ h[h_\ell])$ with $h_k = \langle \dots, \square_k, \dots \rangle$ in the $k^{th}$ position	$\langle q_0, \dots, q_\ell \rangle \mapsto (q_0 \ \dots \ q_\ell)$ $- \mapsto (- \ \dots \ -)$
$\sum_{1 \leq i \leq \ell} K_i(\tau_{i,0}, \dots, \tau_{i,n_i})$	$K_{i_0}$	$(h[h_0] \ \dots \ h[h_{n_{i_0}}])$ with $h_k = K_{i_0}(\dots, \square_k, \dots)$ in the $k^{th}$ position	$K_{i_0}(q_0, \dots, q_{n_{i_0}}) \mapsto (q_0 \ \dots \ q_{n_{i_0}})$ $K_i(\dots) \mapsto \emptyset$ $- \mapsto (- \ \dots \ -)$
$\sum_{1 \leq i \leq \ell} K_i(\tau_{i,0}, \dots, \tau_{i,n_i})$	$\top$	$()$	$K_i(\dots) \mapsto \emptyset$ $- \mapsto ()$

Table 1: Input expansion and pattern specialization –  $\text{EXPAND}^\circ(\tau, tag, h)$ 

the  $\text{COMPILEBRANCH}^\circ$  function.

$$\text{COMPILEBRANCH}^\circ_{\tau_\Delta}(\mathcal{P})(i, tag) = \text{COMPILE}^\circ_{\tau_\Delta}(\mathcal{P}')$$

$$\text{where} \quad \left\{ \begin{array}{l} \text{Let } h_i \text{ the } i\text{-th context in } \mathcal{P} \text{ and } \tau_i = \mathbf{focus}(h_i, \tau_\Delta) \\ \mathcal{H}, f = \text{EXPAND}^\circ(\tau_i, tag, h_i) \\ \mathcal{P}' = \text{FILTERMAPATCOLUMN}(\mathcal{H}, f, \mathcal{P}, i) \end{array} \right.$$

The semantics is given by the function  $\text{EXPAND}^\circ$  defined in Table 1.  $\text{EXPAND}^\circ(\tau, tag, h)$  takes as parameters a tag, the context  $h$  of the considered column and its type  $\tau$ . It returns an expansion  $\mathcal{H}, f$ :  $\mathcal{H}$  defines the new columns and  $f$  defines how to expand the cells of the matrix.  $f$  maps patterns to either vectors or the empty set. Concretely,  $\text{EXPAND}^\circ$  carries out two tasks:

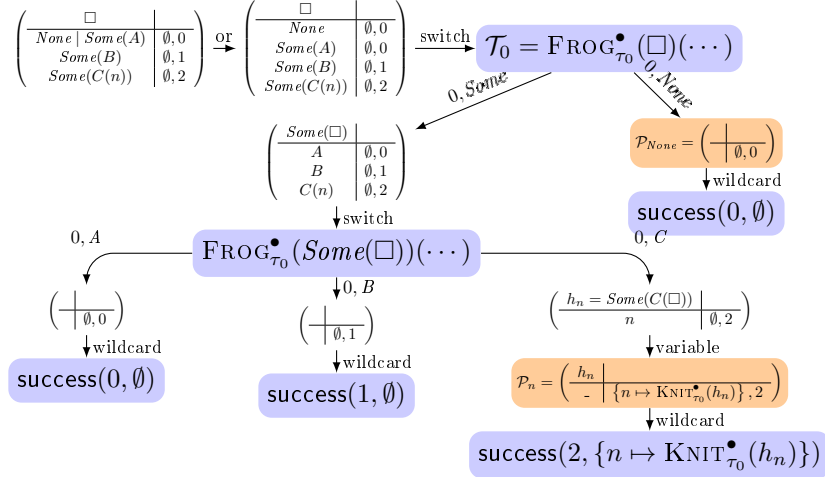
**Specialization** Only keep cases (i.e., rows) in the matrix whose  $i$ -th pattern is compatible with the given tag. For instance, if we consider the branch with tag *Some*, rows with a *None* pattern in position  $i$  are removed. This is done in  $f$  by returning  $\emptyset$  for incompatible cases.

**Expansion** We replace the  $i$ -th column in each input with new columns for the nested patterns. For instance, when considering the tag *Node*, we expand the pattern  $\text{Node}(p_v, p_1, p_2)$ , thus replacing the column  $i$  with three new columns, one for each sub-pattern. This is done by adding new columns in  $\mathcal{H}$  and replacing the column in  $f$ .

The actual new pattern matrix  $\mathcal{P}'$  that serves as input for the recursive call is computed by the operation  $\text{FILTERMAPATCOLUMN}(\mathcal{H}, f, \mathcal{P}, i)$ . It first extends the contexts contained in the column headers by replacing  $h$  with  $\mathcal{H}$ . It then applies  $f$  to each cell  $p_i^j$  of the column  $i$  in matrix  $\mathcal{P}$ . If  $f(p_i^j)$  is the empty set, the whole row at  $j$  is deleted. Otherwise,  $p_i^j$  is replaced with the cells returned by  $f$ . For each cell,  $f$  must return vectors of the length of  $\mathcal{H}$  so that the resulting matrix has a consistent size.

#### 4. Assemble all these sub-trees into a single tree of Switch nodes

We have now gathered all the tags that we need to branch from, along with all the sub-trees corresponding to these branches. We can finally assemble them together. This assembly is representation-specific: indeed, the way to distinguish, for instance, between *Some* and *None* constructors depends on the type under consideration and the choice of representation. We thus delegate the final assembly to the  $\text{FROG}^\bullet$  function.  $\text{FROG}^\bullet_{\tau_\Delta}(h, \mathcal{C})$  takes as parameters the type of the main discriminant  $\tau_\Delta$ , the context  $h$  of the column under consideration, and the list of all cases  $\mathcal{C}$  (i.e., pairs of tags and sub-trees). It returns a tree as defined in Section 4.3.

Figure 13: Pattern compilation for  $m_0 = (None \mid Some(A)) \mid Some(B) \mid Some(C(n))$ .

## 6 End-to-end examples

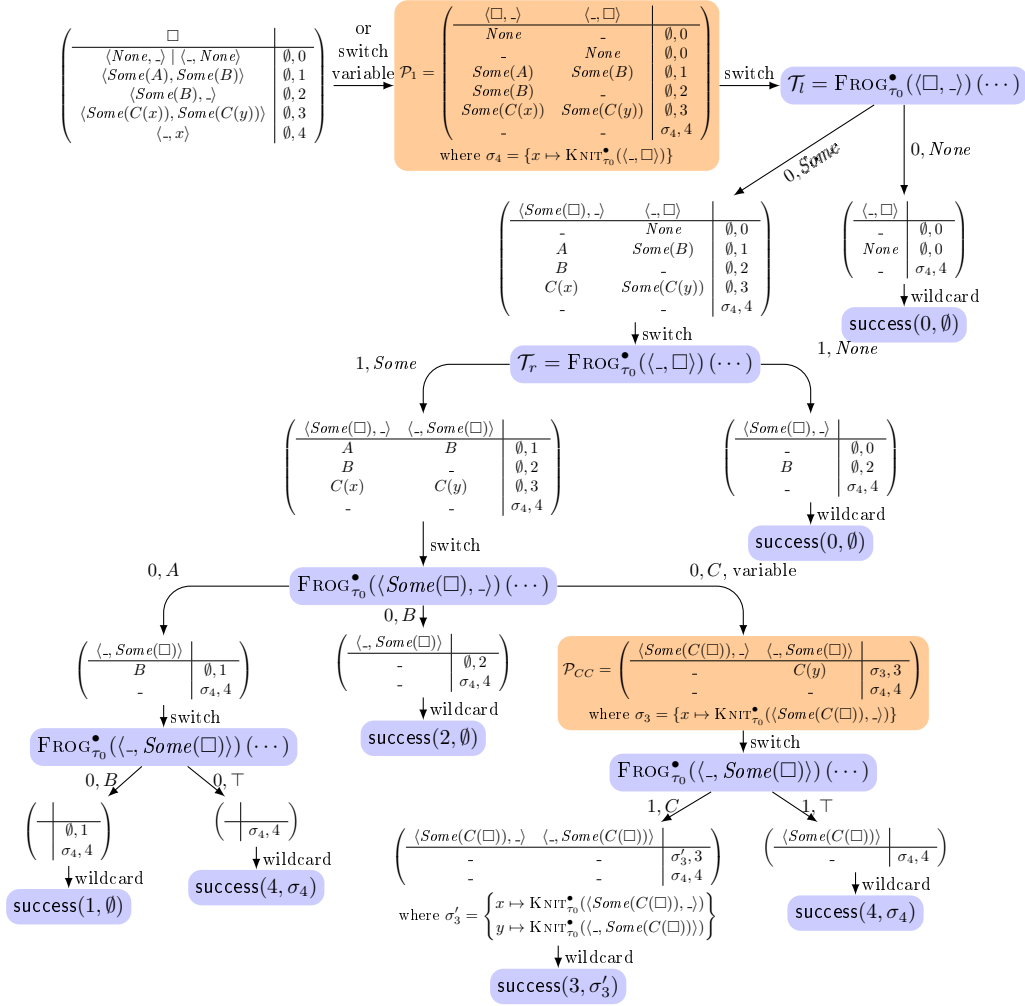
We now define examples of memory representations and how our compilation procedure behaves on them. We first illustrate our compilation procedure on concrete pattern matches. We then showcase a real-world implementation based on the OCaml one. We then define the “packed” representation, which attempts to compact memory values as much as possible.

### 6.1 Compilation procedure

Following our illustration of the boxed representation in Section 4.4, we consider the type  $\tau_0 = None + Some(A + B + C(Int_{32}))$ . The match  $m_0$  and the associated pattern matrix  $\mathcal{P}_0$  are shown below:

$$m_0 = \text{Match} \left\{ \begin{array}{l} | None \mid Some(A) \rightarrow 0 \\ | Some(B) \rightarrow 1 \\ | Some(C(n)) \rightarrow n + 1 \end{array} \right\} \quad \mathcal{P}_0 = \left( \begin{array}{c|c} \square & \\ \hline None \mid Some(A) & \emptyset, 0 \\ Some(B) & \emptyset, 1 \\ Some(C(n)) & \emptyset, 2 \end{array} \right)$$

The derivation of the compilation is presented in Fig. 13. Nodes that result in decision trees are highlighted in blue. Some matrices of particular interest are highlighted in orange and named. The calls to  $\text{FROG}^\bullet$  and  $\text{KNIT}^\bullet$  are left as-is, since we have not specified any memory representation. The algorithm starts with an **or** rule to split the nested patterns into lines, followed by a **switch** rule along the only column of context  $\square$ , thus producing the  $\mathcal{T}_0$  node. The context on which the switch is split is shown in the call to  $\text{FROG}^\bullet$ . The type at this position is of the form  $None + Some(\dots)$ . Furthermore, the *None* and *Some* tags are both present in the column.  $\text{GETTAGS}$  thus gives two branches: *None* and *Some*. The  $\text{COMPILEBRANCH}^\circ$  operation applies input expansion and pattern specialization to obtain the pattern matrices in each of these branches. In the *None* branch, the pattern matrix  $\mathcal{P}_{None}$  has one line (the *None* case in the original code) and no column. Indeed, pattern specialization mandates that only the lines that could match *None* are still present, and there is only one. Input expansion mandates only columns for “interesting” sub-patterns should remain, and there are none. Since all (zero) potential columns are wildcards, we use the **wildcard** rule to produce a **success()** leaf. In the *Some* branch, however, pattern specialization gives us three potential patterns to inspect but removes the outer *Some*, which is now unneeded. The rest of the compilation proceeds by successive application of the **switch** and **wildcard** rules. Finally, in the matrix  $\mathcal{P}_n$ , which corresponds to the context  $h_n = Some(C(\square))$ , we use the **variable** rule to add  $n$  to the binding environment, using  $\text{KNIT}^\bullet$ , before producing the final **success()** leaf.


 Figure 14: Pattern compilation for  $m_1$ .

We now consider a more complex example in Fig. 14 corresponding to the following matching:

$$m_1 = \text{Match} \left\{ \begin{array}{l} | \langle \text{None}, - \rangle | \langle -, \text{None} \rangle \quad \rightarrow \text{None} \\ | \langle \text{Some}(A), \text{Some}(B) \rangle \quad \rightarrow \text{Some}(A) \\ | \langle \text{Some}(B), - \rangle \quad \rightarrow \text{Some}(B) \\ | \langle \text{Some}(C(x)), \text{Some}(C(y)) \rangle \rightarrow \text{Some}(C(x + y)) \\ | \langle -, x \rangle \quad \rightarrow x \end{array} \right\}$$

This example matches on a *tuple* of elements of  $\tau_0$ . In this example, some rule applications were fused for ease of reading. In the first steps, the tuples are translated, thanks to a **switch** case, into multiple columns in the matrix. This is visible in the  $\mathcal{P}_1$  matrix. Each creation of a **switch** node will then pick a given column. For instance, the first  $\mathcal{T}_l$  node picks the first column (context  $\langle \square, - \rangle$ ) and the second  $\mathcal{T}_r$  node picks the second column (context  $\langle -, \square \rangle$ ).

The  $\mathcal{P}_{CC}$  matrix presents a particularity: the first column is composed entirely of wildcards. This particular column can thus not be picked (indeed, the algorithm would loop). We must switch on the second column ( $\langle -, \text{Some}(\square) \rangle$ ). The type at this context is of the form  $A + B + C(\dots)$ ; however, only  $C$  is present in the pattern. GETTAGS thus only produces two branches:  $C$  and  $\top$ .

Many of the matrices where the **wildcard** rule is applied are composed of several lines. This is expected: when looking at the original code, while no case is fully redundant, there are numerous overlaps in which case the first pattern wins. It is thus crucial that the **wildcard** rule only applies to the *first* line of a pattern

matrix.

The rest of the compilation is done by applying the `switch`, `variable` and `wildcard` rules.

## 6.2 The OCaml representation

We now look at a real-world example of memory representation: the OCaml one, depicted in Fig. 15. We picked this representation since it is uniform and well documented Minsky and Madhavapeddy (2021), allowing us to describe it almost completely in our formalism. Indeed, the OCaml representation is similar to the boxed representation described in Section 4.4, in that memory values are either unboxed immediate objects or pointers to blocks, each one-word wide. As this representation is rather uniform and composable, both knitting and frogging remain relatively straightforward.

### 6.2.1 Memory values

As in the boxed representation, OCaml blocks are contiguous word-aligned arrays. Every block starts with a *header* word, whose last byte (called *tag*) provides information about the variant and type of the underlying value. The rest of the header word is composed of information such as block length or GC markers. Since this information is irrelevant for our purposes here (curious readers can enjoy the full details in Minsky and Madhavapeddy (2021)) we will omit them and consider that the header is an integer indicating the variant under consideration. A major difference from the boxed representation is that argument-less variants of sum types are represented by an unboxed integer, rather than empty blocks. The integer is the index of the variant amongst argument-less variants. For instance, *None* will be mapped to the same representation as 0, without any boxing. As a consequence, the type associated with a memory value is no longer sufficient to determine whether it is an unboxed constant or a pointer. In order to distinguish pointers from unboxed values at runtime, the least significant bit of every memory value is used as a *pointer flag*. Concretely, a memory value whose lowest bit is set to 1 contains an unboxed constant in its remaining 63 bits, while a lowest bit set to 0 indicates that the memory value is a block pointer. Blocks are always aligned on 8 bytes, guaranteeing their last bit is always 0. A further consequence is the loss of one bit of precision for constants. This is not problematic for most primitive types<sup>5</sup>. Finally, even though OCaml does offer a concept of “reference”, they are defined as a record (i.e., product types) with a mutable field, and follow the associated representation rather than a separate, pointer-like structure. Reference values are therefore omitted here.

**Example 11** (OCaml representation, memory values). As in Section 4.4, we demonstrate the representation on the type  $\tau_0 = \text{None} + \text{Some}(A + B + C(\text{Int}_{32}))$ . Here are a few examples of values:

$$\begin{aligned} \text{REPR}_{\tau_0}^{\bullet}(\text{None}) &= 0x1 & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(C(n))) &= \&_{64}[0x0; \&_{64}[0x0; ((n \ll 1) \mid 0x1)]]_{64} \\ \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(A)) &= \&_{64}[0x0; 0x01]_{64} & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(B)) &= \&_{64}[0x0; 0x11]_{64} \end{aligned}$$

### 6.2.2 Representation (Knit&Frog)

As with the boxed representation, we first define an auxiliary function  $f$  that inductively deconstructs a context, then use  $f$  to define  $\text{KNIT}^{\bullet}$  and the discriminant portion of  $\text{FROG}^{\bullet}$ .  $f$  is identical to its boxed representation analog. The small optimization on argument-less variants has no effect on knitting, as no context can ever involve such variants (they are, indeed, argument-less).

Frogging undergoes more significant changes. First, numeric constants are transformed to make them match their corresponding memory values with an appropriate shift. Second, since the memory values representing constructor values are no longer uniform across variants, frogging for a sum type potentially yields a two-tiered switch structure. We first consider the case where the considered variants are either all without arguments or all with arguments. The no-arguments case is similar to frogging for primitive types. The case with argument variants involves accessing the lowest byte of the block header after dereferencing the block pointer, then comparing it to the considered constructor indexes. In the general case where both variants with or without arguments are present, the generated code first inspects the lowest bit to determine

<sup>5</sup>Floating-point numbers take up a whole word and are thus boxed in a block whose second word contains the raw, full-width constant.

whether the frogged memory value represents a constructor value with or without arguments; it then branches to a second switch that specifically discriminates between different constructors. The default tag, if present, is propagated to both sub-trees since it accepts either type of variants. Note that this might leave the inner switches with only a single branch (consider, for instance, the option type). For ease of reading, our definition does not remove such single-branch switches, but it is trivial to do so.

**Example 12** (OCaml representation, cont'). As said before, the  $\text{KNIT}^\bullet$  function behaves similarly to the one from the boxed representation. For instance, at the context  $\text{Some}(\square)$ , it dereferences the pointer to access the block, and accesses its first field:  $\text{KNIT}_{\tau_0}^\bullet(\text{Some}(\square)) = *\Delta.2$

$\text{FROG}^\bullet$ , on the other hand, must take into account the particular representation of argument-less variants. For instance, still at position  $\text{Some}(\square)$  with the type  $\tau_0$ , we consider the pattern where the branches  $B$ ,  $C$  and  $\top$  are present. In that case, two switches are required: first to know if the value under consideration is an argument less variant or not by switching on the expression  $*\Delta.2 \ \& \ 0x1$ . If it is not argument-less, it must necessarily be the  $C$  branch. Otherwise, we must test if the variant under consideration is  $B$  or not by switching on  $*\Delta.2$ .

$$\text{FROG}_{\tau_0}^\bullet(\text{Some}(\square))((B, \mathcal{T}_B) \cdot (C, \mathcal{T}_C) \cdot (\top, \mathcal{T}_*)) = \text{switch}(*\Delta.2 \ \& \ 0x1) \begin{cases} 0 \mapsto \mathcal{T}_C \\ 1 \mapsto \text{switch}(*\Delta.2) \\ \quad \begin{cases} 0x11 \mapsto \mathcal{T}_B \\ \top \mapsto \mathcal{T}_* \end{cases} \end{cases}$$

Since the OCaml representation is fairly close to the boxed one, we can state:

**Lemma 2.** The OCaml representation is a valid representation.

### 6.2.3 Compilation

Finally, we present the full compilation of the second example of Section 6.1. The compilation procedure itself was illustrated in Fig. 14 through a graph of the recursive calls of  $\text{COMPILE}^\circ$  where the calls to  $\text{FROG}^\bullet$  and  $\text{KNIT}^\bullet$  were left unresolved. Fig. 16 shows the final decision tree output by our prototype tool, **ribbit**. Rectangle white nodes are  $\text{switch}()$  nodes, with the cases in squares on the right. The green nodes represent  $\text{success}()$  leaves, with the left part being the label, and the right part the binding environment. The choice of order in which to split the columns is identical to the one in Fig. 16. We first inspect each part of the tuple to distinguish between *None* and *Some* (recall that a tuple is a block starting with a tag, hence the members of a couple are at index 1 and 2). In the original derivation, we now need to distinguish between  $A$ ,  $B$  and  $C(\dots)$ . This is done by two switches: first “ $(*(\Delta).1).1 \ \& \ 0x1$ ” to distinguish between argument-less variants  $A + B$  on one hand, and variants with arguments  $C$  on the other; then “ $(*(\Delta).1).1$ ” to distinguish between  $A$  and  $B$ . The rest of the tree proceeds similarly. As before, our procedure doesn't do any sort of common sub-expression elimination, hence a certain redundancy in the result. We assume the follow-up compilation process, which will also choose how to effectively compile the switches, is a better place to apply such optimizations.

## 6.3 The Packed Representation

The Packed Representation is a memory representation which attempts to maximally pack each type, regardless of how it is used: for each type, we define its necessary bit-word length and a function from values to bit-words of that length.

### 6.3.1 Memory values

More formally, let  $\tau$  a type. Let  $\mathcal{V}_\tau = \{v \mid \vdash v : \tau\}$  its values. Since all our base types are finite and there are no recursive types,  $\mathcal{V}_\tau$  is finite: let  $|\mathcal{V}_\tau|$  its cardinal. By definition, there exists an injection, denoted  $\text{REPR}_\tau^\bullet$ , from  $\mathcal{V}_\tau$  to words of size up to  $\lceil \log_2(|\mathcal{V}_\tau|) \rceil$ .

Note that this function is certainly not defined inductively over  $\tau$ ! Additionally, such representation admits no sharing in general, and might have inefficient building and matching of values.

**Example 13** (Packed representation, memory values). For instance, given type  $\tau_0 = \text{None} + \text{Some}(A + B + C(\text{Int}_{32}))$ , the cardinal is  $2^{32} + 3$ , leading to values in 34-bit words, and the following definition:

$$\begin{aligned} \text{REPR}_{\tau_0}^{\bullet}(\text{None}) &= 0 & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(A)) &= 1 \\ \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(B)) &= 2 & \text{REPR}_{\tau_0}^{\bullet}(\text{Some}(C(i))) &= i + 3 \end{aligned}$$

### 6.3.2 Representation (Knit&Frog)

Let us first remark that given a type  $\tau$ , a context  $h$  generates a subset of the image of  $\text{REPR}_{\tau}^{\bullet}$ . Let us denote  $\text{Img}_{\tau}(h)$  this subset.

$\text{KNIT}^{\bullet}$  can then be defined by using the function  $m \mapsto \text{REPR}_{\text{focus}(h,\tau)}^{\bullet}(\text{focus}(h, \text{REPR}_{\tau}^{\bullet-1}(m)))$ , which forms a total mapping from  $\text{Img}_{\tau}(h)$  to  $\text{Img}_{\text{focus}(h,\tau)}(\square)$ . This mapping can be compiled to target expressions using arithmetic and binary operations.

**Example 14** (Packed representation,  $\text{KNIT}^{\bullet}$ ). Following the previous example, the context  $\text{Some}(C(\square))$  corresponds to the subset  $[3 \dots 2^{34} - 1]$ . To map this interval to  $[0 \dots 2^{32} - 1]$ , we can use arithmetic and shifting:  $\text{KNIT}_{\tau_0}^{\bullet}(h, x) = (x - 3) \gg 2$

To build  $\text{FROG}^{\bullet}$ , let us consider a context  $h$ , a list of branches  $l = (t_1, \mathcal{T}_1) \dots (t_n, \mathcal{T}_n)$ , and the patterns  $p_i = \text{Tag2Pat}(t_i)$  built using the function from Definition 1. Since  $\text{REPR}_{\tau}^{\bullet}$  is injective we can use the  $p_i$  patterns to build a corresponding covering partition of  $\text{Img}_{\tau}(h)$  which can then be translated to a decision tree using arithmetic or bit-word tests, whose leaves are the  $\mathcal{T}_i$ .

**Example 15** (Packed representation, cont'). Concretely, here are two examples of patterns following the preceding example:

$$\begin{aligned} \text{FROG}_{\tau_0}^{\bullet}(\square, x)((\text{None}, \mathcal{T}_{\text{None}}), (\text{Some}, \mathcal{T}_{\text{Some}})) &= \text{switch}(x) \begin{cases} 0 \rightarrow \mathcal{T}_{\text{None}} \\ \top \rightarrow \mathcal{T}_{\text{Some}} \end{cases} \\ \text{FROG}_{\tau_0}^{\bullet}(\text{Some}(\square), x)((A, \mathcal{T}_A), (B, \mathcal{T}_B), (C, \mathcal{T}_C)) &= \text{switch}(x) \begin{cases} 0 \rightarrow \text{unreachable} \\ 1 \rightarrow \mathcal{T}_A \\ 2 \rightarrow \mathcal{T}_B \\ \top \rightarrow \mathcal{T}_C \end{cases} \end{aligned}$$

Note how the second switch handles the case 0: indeed, 0, which represents  $\text{None}$ , is not in  $\text{Img}_{\tau_0}(\text{Some}(\square))$  and can thus never be reached.

**Lemma 3.** The packed representation is a valid representation.

As said before, the packed representation defined here is an extreme example, forgoing many optimizations available in more reasonable representations such as sharing or easy subterm extraction. Nevertheless, it shows that our framework is amenable to realistic representation with *some* packing and, given an appropriate definition of  $\text{FROG}^{\bullet}$ , will yield efficient pattern matching code.

## 7 Correctness

We now state the end-to-end correctness of our compilation scheme. The complete proofs are given in Appendix C. Let us first present our hypothesis.

**Exhaustivity** As stated in Section 3, we only consider patterns which are *exhaustive*: i.e. patterns handle all their possible input values.

More formally, let  $p$  a pattern,  $\tau$  a type,  $v$  a value and  $\Gamma$  a typing environment.

$$\vdash p : \tau \rightarrow \Gamma \wedge \vdash v : \tau \implies \exists \sigma. p \triangleright v \rightarrow \sigma \quad (\text{Exhaustivity})$$

Sound analysis for exhaustiveness have already been proposed Maranget (2007), even in the presence of powerful pattern languages Karachalias et al. (2015); Garrigue and Normand (2015). We assume such analysis was used to complete non-exhaustive patterns with wildcards as necessary.



## 7.1 Typing preservation for pattern matching

To ensure that the semantic of pattern is sound with respect to its typing, we state two properties, corresponding to preservation in expression languages.

In order to do so, we first define a typing judgment on values, denoted  $\vdash v : \tau$ , which is defined by a straightforward restriction of the typing on patterns to values. We also equip ourselves with a typing judgment on environments, denoted  $\Gamma \vdash \sigma$ . Both judgments are defined fully in Appendix C.

### Theorem 1. Typing preservation

When a value matches, the bound value environment is well typed with respect to the bound typing environment. More formally, let  $p$  a pattern,  $\tau$  a type,  $v$  a value. Let  $\Gamma$  and  $\sigma$  such that

$$\vdash p : \tau \rightarrow \Gamma \qquad \vdash v : \tau \qquad p \triangleright v \rightarrow \sigma$$

Then, the dynamic and static bound variables coincide:  $\Gamma \vdash \sigma$ .

*Proof.* By induction. See Appendix C.2.

## 7.2 Soundness

### Theorem 2. Equivalence between pattern matching and compilation

Let  $p^1, \dots, p^m$  be  $m$  patterns and  $\tau_\Delta, \gamma$  such that  $\forall j, \vdash p^j : \tau_\Delta \rightarrow \gamma$ . Let  $v$  be a value of type  $\tau_\Delta$ . Let

$$\mathcal{T} = \text{COMPILE}_{\tau_\Delta}^\circ \left( \begin{array}{c|c} \square & \\ \hline p^1 & 1, \emptyset \\ \vdots & \vdots \\ p^m & m, \emptyset \end{array} \right). \text{ Then for any case } j \text{ and store } \sigma:$$

$$p^1 \mid \dots \mid p^m \triangleright v \rightarrow j, \sigma \iff \mathcal{T}[\text{REPR}_{\tau_\Delta}^\bullet(v)] \hookrightarrow j, \{x \mapsto \text{REPR}^\bullet(y)\}_{(x \mapsto y) \in \sigma}$$

*Proof.* Proof sketch. See Appendix C.3 for the full proof.

- We define a *sequence* of pattern matrices from which  $\text{COMPILE}_{\tau_\Delta}^\circ$  generates the nodes that are traversed during execution on input  $\text{REPR}_{\tau_\Delta}^\bullet(v)$ .
- $p^j \triangleright v$  holds iff. any matrix in the sequence contains at least one row of  $\mathcal{P}_k$  that outputs  $j$  and accepts the values obtained by focusing  $v$  on the contexts in  $\mathcal{P}_k$ .
- If the algorithm ends with  $\text{success}(j, s)$  and  $p^j \triangleright v \rightarrow \sigma$ , then  $s$  and  $\sigma$  are compatible.
- The algorithm terminates (i.e., the sequence is finite).
- By induction from the last matrix in the sequence, and using the previous results, the theorem holds for the whole generated decision tree.

□

## 8 Related Works

### 8.1 Memory representation and Algebraic Data Types

Memory representation in functional polymorphic garbage-collected languages was identified quickly as an important area for performance improvements. Peterson (1989) proposes techniques to avoid tagging, while Leroy (1992); Jones and Launchbury (1991) suggest ways to unbox values. Our work encourages new development in this area, by allowing to combine these works with efficient pattern matching compilation. Leroy (1990) presents a calculus which can mix a uniform polymorphic representation and monomorphic optimized representation, which we could use to make several representations cohabit. Colin et al. (2018) details how to extend our source language to handle recursive types in the presence of unboxing. Many of these works are implemented in some capacity in OCaml and Haskell.

Iannetta et al. (2021); Koparkar et al. (2021) propose drastically different representations for Algebraic Data Types, where almost everything is flattened, allowing excellent cache behavior and parallelism. Our work would augment these approaches with efficient decision tree generation.

## 8.2 Pattern languages

Pattern languages for algebraic data types were first introduced by the HOPE language Burstall et al. (1980). Its general form has been adopted mostly as-is in mainstream languages with rich static typing such as Haskell, OCaml, F#, Scala or Rust, but also more recently in more general languages such as Python and soon Java. This diversity of host languages, with their very varied compilation techniques and memory representations, make our framework all the more relevant.

We focused on the core of pattern matching language, with some minor extensions like disjunctive patterns. Other extensions include ranges, guards, matching of polymorphic variants Garrigue (1998), and exception patterns (in recent OCaml versions). These are orthogonal to our work.

Pattern matching is also used pervasively in dependently typed languages Cockx et al. (2016); Tuerk et al. (2015). Since our definition of  $\text{KNIT}^\bullet$  and  $\text{FROG}^\bullet$  takes the type into account, we can't simply delegate the typing to a previous phase and operate on an untyped representation. One difficulty in dependent pattern matching (and in GADTs Garrigue and Normand (2015); Karachalias et al. (2015)) is that matching on a column can reveal the type of other columns. To adapt our framework to this setting, we could keep track of columns whose type is not revealed yet, and prioritize the decomposition of other columns. Columns would still have a coherent type, thanks to the split of the pattern matrix. OCaml implements a simplified version of this (it always splits on the first column, and type dependency can only go towards previous columns).

Active patterns Syme et al. (2007) from F# allow users to abstract over patterns by exposing “constructors” which do not directly reflect the underlying definition of the algebraic data type. The original formulation allows active patterns to execute arbitrary code during matching, making optimization difficult. Pattern synonyms Pickering et al. (2016) restrict patterns to be closer to the actual implementation, making it more amenable to optimization procedures such as ours.

## 8.3 Pattern Matching Optimization

Optimization of pattern matching for algebraic data types was first proposed by Augustsson (1985) in the context of the LML language. It first introduced the notion of pattern matrices and produced a backtracking automaton. Optimization of pattern matching is significantly more delicate in *lazy* languages (such as LML), since matching should avoid forcing unnecessary terms. Several semantics and associated optimizations have been proposed (Puel and Suárez, 1993; Maranget, 1992; Wadler, 1987), the last of which is used the current implementation of GHC (Marlow et al., 2004).

For *strict* languages however, Baudinet and MacQueen (1985) remark that the different parts of the discriminant can be evaluated in any order, allowing for more optimization opportunities. Fessant and Maranget (2001) first proposed optimizations for backtracking automata, introducing the “row and column” approach to split the pattern matrix. Their technique is currently used in OCaml. This approach was later refined by Maranget (2008) to produce good decision trees, which we base our work on. It delivers excellent performance, while being reasonable to compute in practice. Sestoft (1996) emits a rough tree of `if` nodes, and relies on a global supercompilation pass to optimize the resulting tree. It is not clear how to make it parametric in term of the memory representation. Kosarev et al. (2020) explore a different optimization technique by encoding the choice of optimal decision tree into a relational synthesis problem, and solving through miniKanren. Their idea is very promising, but fails to scale to big matches. Solodkyy et al. (2013) propose patterns-as-library for C++ based on objects and template meta-programming. In all these cases, DAGs with maximal sharing can be created from trees by using Hash-consing Filliâtre and Conchon (2006).

Most approaches rely on heuristics for the choice of column to split (albeit with a limited choice for lazy languages). Maranget (2008) introduces the “Necessity” heuristic. A study of heuristics is done in Scott and Ramsey (2000). Both conclude that the choice of heuristic only has minor performance consequences in most cases, but can matter for very particular matches. We believe the choice of memory representation is a much bigger factor.

## 9 Conclusion

We have presented Knit&Frog, a technique to compile pattern matching parameterized by an arbitrary memory representation. On our way, we have given a novel way to define memory representations, along

with a validity criterion that ensures such a representation is appropriate for pattern matching. We used this criterion to prove the correctness of our compilation procedure. Finally, we illustrated our approach on several examples, including a realistic description of the OCaml memory representation. To our knowledge, this is the first generic description of memory representation, and the relation with Algebraic Data Types. We have also implemented our technique in a prototype tool called `ribbit` and shown its output on concrete examples.

Our technique paves the way towards the formalization and description of new optimization techniques for memory representation. In recent versions, Rust has been introducing more and more complex memory representation optimizations, which are so far unspecified. Furthermore, we believe there is a trove of untapped optimization yet to be explored when it comes to the memory representation of Algebraic Data Types. Some promising leads would be to apply super-optimization to individual performance-sensitive data structures, or to allow programmers to specify whether types should be optimized for space, cache behavior, or even best sharing. We hope this work serves as a stepping stone for these further optimizations.

$$\text{Tag2Word}_\tau(c) = (c \ll 1 \mid 0\text{x}1) \quad \text{when } \tau = T \in \mathcal{C}$$

$$\text{Tag2Word}_{\sum K_i(\overline{\tau_{i,j}}^{J(i)})}(K_{i_0}) = \begin{cases} \text{Tag2Word}_{\text{Int}}(\text{nth}(K, \{K_i \mid J(i) = 0\})) & \text{when } J(i_0) = 0 \\ \text{nth}(K, \{K_i \mid J(i) > 0\}) & \text{when } J(i_0) > 0 \end{cases}$$

(a) The Tag2Word utility function

$$\text{REPR}_T^\bullet(c) = \text{Tag2Word}_T(c) \quad \text{REPR}_{\prod \overline{\tau_i}}^\bullet(\overline{\langle v_i \rangle}) = \&_{64} \left[ 0; \overline{\text{REPR}_{\tau_i}^\bullet(v_i)} \right]_{64}$$

$$\text{REPR}_\tau^\bullet(K_i) = \text{Tag2Word}_\tau(K_i) \quad \text{REPR}_{\sum K_i(\overline{\tau_{i,j}})}^\bullet(K_i(\overline{v_j})) = \&_{64} \left[ \text{Tag2Word}_\tau(K_i); \overline{\text{REPR}_{\tau_{i,j}}^\bullet(v_j)} \right]_{64}$$

(b) Memory values

$$f(\square, e) = e \quad f(\langle -, \dots, h_k, \dots, - \rangle, e) = f(K(\langle -, \dots, h_k, \dots, - \rangle), e) = f(h_k, k + 1. * e)$$

$$\text{KNIT}_\tau^\bullet(h) = f(h, \Delta)$$

(c) Knitting

$$\text{FROG}_\tau^\bullet(h) (\overline{t_k}, \overline{\mathcal{T}_k}) = \text{switch}(f(h, \Delta)) \left\{ \overline{\text{Tag2Word}_\tau(t_k)} \mapsto \overline{\mathcal{T}_k} \right.$$

where **focus**  $(h, \tau) = T \in \mathcal{C}$

$$\text{FROG}_\tau^\bullet(h) (\overline{(K_{i_k}, \mathcal{T}_{i_k})} \cdot (\top, \mathcal{T}_*)) = \text{switch}(f(h, \Delta)) \left\{ \overline{\text{Tag2Word}_\tau(K_{i_k})} \mapsto \overline{\mathcal{T}_{i_k}} \right.$$

where **focus**  $(h, \tau) = \sum K_i(\overline{\tau_{i,j}}^{J(i)})$  and  $\forall k, J(i_k) = 0$

$$\text{FROG}_\tau^\bullet(h) (\overline{(K_{i_k}, \mathcal{T}_{i_k})} \cdot (\top, \mathcal{T}_*)) = \text{switch}(f(h, \Delta) \& 0\text{x}1) \begin{cases} 1 \mapsto \text{FROG}_\tau^\bullet(h) (\overline{(K_{i_k}, \mathcal{T}_{i_k})}_{J(i_k)=0} \cdot (\top, \mathcal{T}_*)) \\ 0 \mapsto \text{switch}(* (f(h, \Delta)).1) \begin{cases} \overline{\text{Tag2Word}_\tau(K_{i_k})} \mapsto \overline{\mathcal{T}_{i_k}} \\ \top \mapsto \mathcal{T}_* \end{cases} \end{cases}$$

where **focus**  $(h, \tau) = \sum K_i(\overline{\tau_{i,j}}^{J(i)})$

(d) Frogging

Figure 15: The OCaml representation

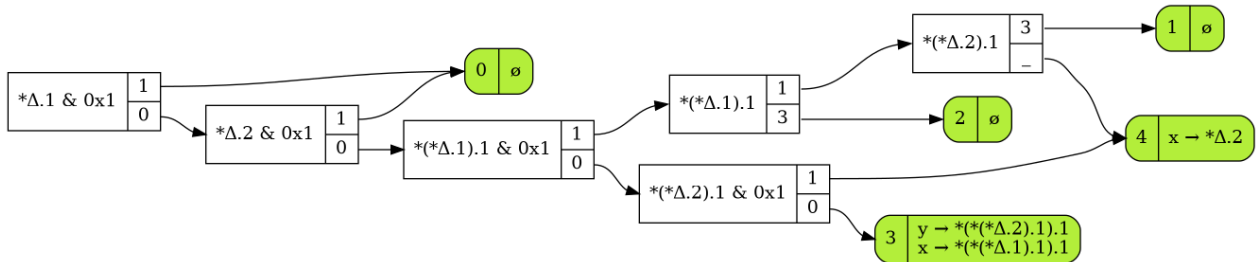


Figure 16: The decision tree for Fig. 14 with the OCaml representation, using our prototype tool *ribbit*.

## References

- Lennart Augustsson. 1985. Compiling Pattern Matching. In *Functional Programming Languages and Computer Architecture, FPCA 1985, Nancy, France, September 16-19, 1985, Proceedings (Lecture Notes in Computer Science, Vol. 201)*, Jean-Pierre Jouannaud (Ed.). Springer, 368–381. [https://doi.org/10.1007/3-540-15975-4\\_48](https://doi.org/10.1007/3-540-15975-4_48)
- Marianne Baudinet and David MacQueen. 1985. Tree pattern matching for ML. (1985).
- Rod M. Burstall, David B. MacQueen, and Donald Sannella. 1980. HOPE: An Experimental Applicative Language. In *Proceedings of the 1980 LISP Conference, Stanford, California, USA, August 25-27, 1980*. ACM, 136–143. <https://doi.org/10.1145/800087.802799>
- Jesper Cockx, Dominique Devriese, and Frank Piessens. 2016. Eliminating dependent pattern matching without K. *J. Funct. Program.* 26 (2016), e16. <https://doi.org/10.1017/S0956796816000174>
- Simon Colin, Rodolphe Lepigre, and Gabriel Scherer. 2018. Unboxing Mutually Recursive Type Definitions in OCaml. *arXiv preprint arXiv:1811.02300* (2018).
- Fabrice Le Fessant and Luc Maranget. 2001. Optimizing Pattern Matching. In *Proceedings of the Sixth ACM SIGPLAN International Conference on Functional Programming (ICFP '01), Firenze (Florence), Italy, September 3-5, 2001*, Benjamin C. Pierce (Ed.). ACM, 26–37. <https://doi.org/10.1145/507635.507641>
- Jean-Christophe Filliâtre and Sylvain Conchon. 2006. Type-safe modular hash-consing. In *Proceedings of the ACM Workshop on ML, 2006, Portland, Oregon, USA, September 16, 2006*, Andrew Kennedy and François Pottier (Eds.). ACM, 12–19. <https://doi.org/10.1145/1159876.1159880>
- Jacques Garrigue. 1998. Programming with polymorphic variants. In *ML Workshop*, Vol. 13. Baltimore.
- Jacques Garrigue and Jacques Le Normand. 2015. GADTs and Exhaustiveness: Looking for the Impossible. In *Proceedings ML Family / OCaml Users and Developers workshops, ML Family/OCaml 2015, Vancouver, Canada, 3rd & 4th September 2015 (EPTCS, Vol. 241)*, Jeremy Yallop and Damien Doligez (Eds.). 23–35. <https://doi.org/10.4204/EPTCS.241.2>
- Paul Iannetta, Laure Gonnord, and Gabriel Radanne. 2021. Compiling pattern matching to in-place modifications. In *GPCE '21: Concepts and Experiences, Chicago, IL, USA, October 17 - 18, 2021*, Eli Tilevich and Coen De Roover (Eds.). ACM, 123–129. <https://doi.org/10.1145/3486609.3487204>
- Simon L. Peyton Jones and John Launchbury. 1991. Unboxed Values as First Class Citizens in a Non-Strict Functional Language. In *Functional Programming Languages and Computer Architecture, 5th ACM Conference, Cambridge, MA, USA, August 26-30, 1991, Proceedings (Lecture Notes in Computer Science, Vol. 523)*, John Hughes (Ed.). Springer, 636–666. [https://doi.org/10.1007/3540543961\\_30](https://doi.org/10.1007/3540543961_30)
- Georgios Karachalias, Tom Schrijvers, Dimitrios Vytiniotis, and Simon L. Peyton Jones. 2015. GADTs meet their match: pattern-matching warnings that account for GADTs, guards, and laziness. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, ICFP 2015, Vancouver, BC, Canada, September 1-3, 2015*, Kathleen Fisher and John H. Reppy (Eds.). ACM, 424–436. <https://doi.org/10.1145/2784731.2784748>
- Chaitanya Koparkar, Mike Rainey, Michael Vollmer, Milind Kulkarni, and Ryan R. Newton. 2021. Efficient Tree-Traversals: Reconciling Parallelism and Dense Data Representations. *Proc. ACM Program. Lang.* 5, ICFP, Article 91 (aug 2021), 29 pages. <https://doi.org/10.1145/3473596>
- Dmitry Kosarev, Petr Lozov, and Dmitry Boulytchev. 2020. Relational Synthesis for Pattern Matching. In *Processings of the 18th Programming Languages and Systems Asian Symposium, APLAS 2020, Fukuoka, Japan, 2020 (Lecture Notes in Computer Science, Vol. 12470)*, Bruno C. d. S. Oliveira (Ed.). Springer, 293–310. [https://doi.org/10.1007/978-3-030-64437-6\\_15](https://doi.org/10.1007/978-3-030-64437-6_15)

- Xavier Leroy. 1990. Efficient Data Representation in Polymorphic Languages. In *Programming Language Implementation and Logic Programming, 2nd International Workshop PLILP'90, Linköping, Sweden, August 20-22, 1990, Proceedings (Lecture Notes in Computer Science, Vol. 456)*, Pierre Deransart and Jan Maluszynski (Eds.). Springer, 255–276. <https://doi.org/10.1007/BFb0024189>
- Xavier Leroy. 1992. Unboxed Objects and Polymorphic Typing. In *Conference Record of the Nineteenth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Albuquerque, New Mexico, USA, January 19-22, 1992*, Ravi Sethi (Ed.). ACM Press, 177–188. <https://doi.org/10.1145/143165.143205>
- Luc Maranget. 1992. Compiling Lazy Pattern Matching. In *Proceedings of the Conference on Lisp and Functional Programming, LFP 1992, San Francisco, California, USA, 22-24 June 1992*, Jon L. White (Ed.). ACM, 21–31. <https://doi.org/10.1145/141471.141499>
- Luc Maranget. 2007. Warnings for pattern matching. *J. Funct. Program.* 17, 3 (2007), 387–421. <https://doi.org/10.1017/S0956796807006223>
- Luc Maranget. 2008. Compiling pattern matching to good decision trees. In *Proceedings of the ACM Workshop on ML, 2008, Victoria, BC, Canada, September 21, 2008*, Eijiro Sumii (Ed.). ACM, 35–46. <https://doi.org/10.1145/1411304.1411311>
- Simon Marlow, Simon Peyton Jones, et al. 2004. The glasgow haskell compiler.
- Yaron Minsky and Anil Madhavapeddy. 2021. *Real World OCaml*. Chapter Memory Representation of Values. <https://dev.realworldocaml.org/runtime-memory-layout.html>
- John Peterson. 1989. Untagged Data in Tagged Environments: Choosing Optimal Representations at Compile Time. In *Proceedings of the fourth international conference on Functional programming languages and computer architecture, FPCA 1989, London, UK, September 11-13, 1989*, Joseph E. Stoy (Ed.). ACM, 89–99. <https://doi.org/10.1145/99370.99377>
- Matthew Pickering, Gergo Érdi, Simon Peyton Jones, and Richard A. Eisenberg. 2016. Pattern synonyms. In *Proceedings of the 9th International Symposium on Haskell, Haskell 2016, Nara, Japan, September 22-23, 2016*, Geoffrey Mainland (Ed.). ACM, 80–91. <https://doi.org/10.1145/2976002.2976013>
- Laurence Puel and Ascánder Suárez. 1993. Compiling Pattern Matching by Term Decomposition. *J. Symb. Comput.* 15, 1 (1993), 1–26. <https://doi.org/10.1006/jscs.1993.1001>
- Kevin Scott and Norman Ramsey. 2000. When do match-compilation heuristics matter. *University of Virginia, Charlottesville, VA* (2000).
- Peter Sestoft. 1996. ML pattern match compilation and partial evaluation. In *Partial Evaluation*. Springer, 446–464.
- Yuriy Solodkyy, Gabriel Dos Reis, and Bjarne Stroustrup. 2013. Open pattern matching for C++. In *Generative Programming: Concepts and Experiences, GPCE'13, Indianapolis, IN, USA - October 27 - 28, 2013*, Jaakko Järvi and Christian Kästner (Eds.). ACM, 33–42. <https://doi.org/10.1145/2517208.2517222>
- Don Syme, Gregory Neverov, and James Margetson. 2007. Extensible pattern matching via a lightweight language extension. In *Proceedings of the 12th ACM SIGPLAN International Conference on Functional Programming, ICFP 2007, Freiburg, Germany, October 1-3, 2007*, Ralf Hinze and Norman Ramsey (Eds.). ACM, 29–40. <https://doi.org/10.1145/1291151.1291159>
- Thomas Tuerk, Magnus O. Myreen, and Ramana Kumar. 2015. Pattern Matches in HOL: - A New Representation and Improved Code Generation. In *Interactive Theorem Proving - 6th International Conference, ITP 2015, Nanjing, China, August 24-27, 2015, Proceedings (Lecture Notes in Computer Science, Vol. 9236)*, Christian Urban and Xingyuan Zhang (Eds.). Springer, 453–468. [https://doi.org/10.1007/978-3-319-22102-1\\_30](https://doi.org/10.1007/978-3-319-22102-1_30)

Philip Wadler. 1987. Efficient compilation of pattern-matching. *The implementation of functional programming languages* (1987).

Pierre Wilke. 2016. *Formally verified compilation of low-level C code. (Compilation formellement vérifiée de code C de bas-niveau)*. Ph.D. Dissertation. University of Rennes 1, France. <https://tel.archives-ouvertes.fr/tel-01483676>

## A Typing of values

In this section, we define the typing judgment  $\vdash v : \tau$  for values and  $\vdash \sigma : \Gamma$  for environments in Fig. 17. The rules are a direct subset of the rules for patterns presented in Section 3.

$$\begin{array}{c}
\text{CONSTANT} \\
\frac{c \in T}{\vdash c : T}
\end{array}
\qquad
\begin{array}{c}
\text{REFERENCE} \\
\frac{\vdash v : \tau}{\vdash \&v : \&\tau}
\end{array}
\qquad
\begin{array}{c}
\text{TUPLE} \\
\frac{\forall i, \vdash v_i : \tau_i}{\vdash \langle \bar{v}_i \rangle : \prod \tau_i}
\end{array}
\qquad
\begin{array}{c}
\text{CONSTRUCTOR} \\
\frac{\exists i, K = K_i \quad \forall j, \vdash v_j : \tau_{i,j}}{\vdash K(\bar{v}_j) : \sum K_i(\tau_i, j)}
\end{array}$$

$$\begin{array}{c}
\text{ENV} \\
\frac{\forall i, \vdash v_i : \tau_i}{\vdash \{x_i \mapsto v_i\} : \{x_i : \tau_i\}}
\end{array}$$

Figure 17: Typing of values and environments –  $\vdash v : \tau$

## B Contexts and focusing

We initially introduced the focus judgment in Section 4.2.2 as a way to manipulate values and types with respect to a context  $h$ . We defined two operations: **focus**  $(h, v)$  which focuses a value and returns a value, and **focus**  $(h, \tau)$  which focuses a type and returns a type. For the purpose of the proof, we introduce an additional one: **focus**  $(h, p)$ , which focuses a pattern and returns a pattern, by collecting all the alternatives. The complete definitions are given in Fig. 18.

## C Proof of correctness

We now give the proof of correctness of the theorems stated in Section 7: the typing preservation, and the soundness of our compilation procedure. We also define a few convenience results on the way.

### C.1 Coherence of focusing

We use focusing pervasively in our proof. Let us state a convenience result which shows that focus and typing of values are coherent:

**Lemma 4.** Let  $h$  a context,  $v$  a value,  $\tau$  a type.

$$\tau \vdash v \wedge \mathbf{focus}(h, v) = v' \wedge \mathbf{focus}(h, \tau) = \tau' \implies \tau' \vdash v'$$

*Proof.* By induction over the typing judgment. □

### C.2 Typing preservation

*Proof of Theorem 1.* Let  $p$  a pattern,  $\tau$  a type,  $v$  a value. Let  $\Gamma$  and  $\sigma$  such that

$$\vdash p : \tau \rightarrow \Gamma \qquad \vdash v : \tau \qquad p \triangleright v \rightarrow \sigma$$

By induction on  $p$ , we show that  $x \in \Gamma \iff x \in \sigma$  and  $\forall x \in \Gamma, \vdash \sigma(x) : \Gamma(x)$ .

**Wildcard or constant pattern** no bound variables:  $\Gamma = \sigma = \emptyset$ .

**Variable pattern** ( $p = x$ )  $x$  is the only bound variable:  $\Gamma = \{x : \tau\}$ ;  $\sigma = \{x \mapsto v\}$ . Since  $\vdash v : \tau$ , we have  $\vdash \sigma(x) : \tau$ .



$$\begin{aligned} \mathbf{focus}(\square, v) &= v & \mathbf{focus}(\&h, \&v) &= \mathbf{focus}(h, v) \\ \mathbf{focus}(\langle -, \dots, -, h, -, \dots, - \rangle, \langle -, \dots, -, v, -, \dots, - \rangle) &= \mathbf{focus}(h, v) \\ \mathbf{focus}(K(-, \dots, -, h, -, \dots, -), K(-, \dots, -, v, -, \dots, -)) &= \mathbf{focus}(h, v) \end{aligned}$$

(a) Focus on values –  $\mathbf{focus}(h, v)$ 

$$\begin{aligned} \mathbf{focus}(\square, p) &= p & \mathbf{focus}(\&h, \&p) &= \mathbf{focus}(h, p) \\ \mathbf{focus}(\langle -, \dots, -, h, -, \dots, - \rangle, \langle -, \dots, -, p, -, \dots, - \rangle) &= \mathbf{focus}(h, p) \\ \mathbf{focus}(K(-, \dots, -, h, -, \dots, -), K(-, \dots, -, p, -, \dots, -)) &= \mathbf{focus}(h, p) \\ \mathbf{focus}(h, p_1 \mid p_2) &= \mathbf{focus}(h, p_1) \mid \mathbf{focus}(h, p_2) \end{aligned}$$

(b) Focus on patterns –  $\mathbf{focus}(h, p)$ 

$$\mathbf{focus}(\square, \tau) = \tau \quad \mathbf{focus}(\&h, \&\tau) = \mathbf{focus}(h, \tau) \quad \mathbf{focus}\left(\left\langle -, \dots, -, h_k, -, \dots, - \right\rangle, \prod \tau_i\right) = \mathbf{focus}(h, \tau_k)$$

$$\mathbf{focus}\left(K_{i_0}(-, \dots, -, h_k, -, \dots, -), \sum_{0 \leq i < n} K_i(\overline{\tau_{i,j}})\right) = \mathbf{focus}(h, \tau_{i_0,k}) \text{ with } 0 \leq i_0 < n$$

(c) Focus on types –  $\mathbf{focus}(h, \tau)$ 

Figure 18: Focusing judgments

**Reference** ( $p = \&p'$ ) let  $v'$  and  $\tau'$  such that  $\vdash p' : \tau', \Gamma$ ,  $v = \&v'$  and  $p' \triangleright v', \sigma$ . From the Lemma 4, we have  $\vdash v' : \tau'$ . Any variable bound in  $p$  is bound in  $p'$ ; from the induction hypothesis, we have  $x \in \Gamma \iff x \in \sigma$  and  $\forall x \in \Gamma, \vdash \sigma(x) : \Gamma(x)$  for any bound variable  $x$ .

**Tuple** ( $p = \langle \overline{p_k} \rangle$ ) let  $\overline{\tau_k}, \overline{v_k}, \overline{\Gamma_k}$  and  $\overline{\sigma_k}$  such that:

$$\forall k, \vdash p_k : \tau_k, \Gamma_k \quad \forall k, p_k \triangleright v_k, \sigma_k \quad v = \langle \overline{v_k} \rangle \quad \forall k, \vdash v_k : \tau_k \text{ (Lemma 4)} \quad \tau = \prod \tau_k \quad \Gamma = \bigsqcup \Gamma_k$$

$$\sigma = \bigsqcup \sigma_k$$

For any bound variable  $x \in \Gamma$ , there is exactly one  $k$  such that  $x \in \Gamma_k$ ,  $\Gamma(x) = \Gamma_k(x)$  and  $\sigma(x) = \sigma_k(x)$ . From the induction hypothesis, we have  $x \in \sigma_k$  and  $\vdash \sigma_k(x) : \Gamma_k(x)$ .

Conversely, for any  $x \in \sigma$ , there is exactly one  $k$  such that  $x \in \sigma_k$ ; from the induction hypothesis, we have  $x \in \Gamma_k \subset \Gamma$ .

**Constructor pattern** ( $p = K(\overline{p_k})$ ) let  $\tau = \sum K_i(\overline{\tau_{i,k}})$  and  $i$  such that  $K = K_i$  and  $\forall k, \vdash p_k : \tau_k$ . Let  $\overline{v_k}, \overline{\Gamma_k}$  and  $\overline{\sigma_k}$  such that:

$$\forall k, v = K(\overline{v_k}) \quad \forall k, \vdash v_k : \tau_k \text{ (Lemma 4)} \quad \forall k, \vdash p_k : \tau_k, \Gamma_k \quad \forall k, p_k \triangleright v_k, \sigma_k \quad \Gamma = \bigcup \Gamma_k \quad \sigma = \bigcup \sigma_k$$

We show that the properties hold as in the previous case.

**Disjunction** ( $p = p_1 \mid p_2$ ) we have  $\vdash p_1 : \tau, \Gamma$ ,  $\vdash p_2 : \tau, \Gamma$  and either  $p_1 \triangleright v, \sigma$  or  $p_2 \triangleright v, \sigma$ . The variables bound in  $p$  are bound in  $p_1$  or  $p_2$ . In both cases, from the induction hypothesis, we have  $x \in \Gamma \iff x \in \sigma$  and  $\forall x \in \Gamma, \vdash \sigma(x) : \Gamma(x)$ .

□

### C.3 Soundness

*Proof of Theorem 2.* Let  $p_0^0, \dots, p_0^{m-1}$  be  $m$  (initial) patterns and  $\tau_\Delta, \gamma$  such that  $\forall j, \vdash p_0^j : \tau_\Delta \rightarrow \gamma$ . Let  $v$

be a value such that  $\vdash v : \tau_\Delta$  and  $r = \text{REPR}_{\tau_\Delta}^\bullet(v)$ . Finally, let  $\mathcal{P}_0 = \left( \begin{array}{c|c} \square & 0, \emptyset \\ \hline p_0^0 & \vdots \\ \vdots & \vdots \\ p_0^{m-1} & m-1, \emptyset \end{array} \right)$ .

We show that for any  $j, \sigma, p_0^0 | \dots | p_0^{m-1} \triangleright v \rightarrow j, \sigma$  iff.  $\text{COMPILE}_{\tau_\Delta}^\circ(\mathcal{P}_0)[r] \hookrightarrow j, \{x \mapsto \text{REPR}^\bullet(y)\}_{x \in \sigma}$ .

#### C.3.1 Compilation steps

We define a sequence of pattern matrices corresponding to the  $\text{COMPILE}_{\tau_\Delta}^\circ$  calls that define those nodes that are part of the execution path on input  $r$ .

For any pattern matrix  $\mathcal{P}_k$  in the sequence, we write  $\mathcal{P}_k = \left( \begin{array}{ccc|c} h_{k,0} & \dots & h_{k,n_k-1} & \\ \hline p_{k,0}^0 & \dots & p_{k,n_k-1}^0 & j_k^0, s_k^0 \\ \vdots & \ddots & \vdots & \vdots \\ p_{k,0}^{m_k-1} & \dots & p_{k,n_k-1}^{m_k-1} & j_k^{m_k-1}, s_k^{m_k-1} \end{array} \right)$ .

We define the sequence of relevant pattern matrices inductively:

- The first pattern matrix is  $\mathcal{P}_0$ , as defined above.
- If  $\mathcal{P}_k$  belongs to the sequence, we follow the definition of  $\text{COMPILE}_{\tau_\Delta}^\circ$  to determine the remaining steps:
  - If  $\mathcal{P}_k = \emptyset / m_k = 0$ , the sequence ends (no pattern case).
  - If  $\mathbf{Pos}_k^0 = \emptyset$ , the sequence ends (wildcard case).
  - If the variable (resp. disjunction) case applies and we lift (resp. split) the pattern in column  $i$  and row  $\ell$ , we define the next matrix  $\mathcal{P}_{k+1}$  according to  $\text{COMPILE}_{\tau_\Delta}^\circ$  and write  $\mathcal{P}_k \rightarrow_{\text{var},i,\ell} \mathcal{P}_{k+1}$  (resp.  $\mathcal{P}_k \rightarrow_{\text{or},i,\ell} \mathcal{P}_{k+1}$ ).
  - Otherwise, the switch case applies. Let  $i$  the column chosen by `PICKCOLUMN`. We have  $h_{k,i} \in \mathbf{Pos}_k^1$ . Let  $\tau_i = \mathbf{focus}(h_{k,i}, \tau_\Delta)$  and  $t$  be the first tag generated by `GETTAGS`[ $\tau_i$ ] that matches  $\mathbf{focus}(h_{k,i}, v)$ . `FROG`<sup>•</sup> (used in the last step of the switch case of  $\text{COMPILE}_{\tau_\Delta}^\circ$ ) guarantees that at this point, the execution path on input  $r$  will branch to the subtree associated with  $t$ . This subtree is generated by  $\text{COMPILEBRANCH}^\circ[\tau_\Delta](\mathcal{P}_k)(i, t)$ , which performs one recursive call to  $\text{COMPILE}_{\tau_\Delta}^\circ$ . Let  $\mathcal{P}_{k+1}$  be the pattern matrix built by  $\text{COMPILEBRANCH}^\circ$  and passed to  $\text{COMPILE}_{\tau_\Delta}^\circ$  for this recursive call. Then  $\mathcal{P}_{k+1}$  is the next matrix in the sequence, and we write  $\mathcal{P}_k \rightarrow_{\text{switch},i,t} \mathcal{P}_{k+1}$ .

Finally, for any matrix  $\mathcal{P}_k$  in the sequence and initial case  $j$ , let  $\mathcal{P}_k^j$  be the (possibly empty) sub-matrix of  $\mathcal{P}_k$  that contains all rows such that  $j_k^\ell = j$ . We trivially show by induction on  $k$  that these sub-matrices

never overlap:  $\mathcal{P}_k = \left( \begin{array}{c} \mathcal{P}_k^0 \\ \vdots \\ \mathcal{P}_k^{m-1} \end{array} \right)$ .

#### C.3.2 Case matching equivalence

**Proposition 1.** *Let  $j$  an initial case and  $\mathcal{P}_k$  a pattern matrix in the sequence. We have  $p_0^j \triangleright v$  iff.  $\exists \ell, j_k^\ell = j \wedge \forall i, p_{k,i}^\ell \triangleright \mathbf{focus}(h_{k,i}, v)$ .*

*Proof.* By induction on  $k$ . The base case is trivial ( $n_0 = 1, p_{0,0}^j = p_0^j$  and  $h_{0,0} = \square$ ). We now consider the case where  $k > 0$  and suppose we have

$$p_0^j \triangleright v \iff \exists \ell, j_{k-1}^\ell = j \wedge \forall i, p_{k-1,i}^\ell \triangleright \mathbf{focus}(h_{k-1,i}, v)$$

We show that the property carries over to the next pattern matrix  $\mathcal{P}_k$  by induction on the transition step from  $\mathcal{P}_{k-1}$  to  $\mathcal{P}_k$ :

**Variable case:**  $\mathcal{P}_{k-1} \rightarrow_{\text{var},i,\ell} \mathcal{P}_k$ . We have  $p_{k,i}^{\ell'} = p_{k-1,i'}^{\ell'}$  if  $i' \neq i$  or  $\ell' \neq \ell$ ,  $p_{k,i}^{\ell} = -$  and  $p_{k-1,i}^{\ell}$  is a variable pattern. In both cases, we have

$$p_{k,i}^{\ell} \triangleright \mathbf{focus}(h_{k,i}, v) \iff p_{k-1,i}^{\ell} \triangleright \mathbf{focus}(h_{k-1,i}, v)$$

therefore, the property holds for  $\mathcal{P}_k$ .

**Disjunction case:**  $\mathcal{P}_{k-1} \rightarrow_{\text{or},i,\ell} \mathcal{P}_k$ . We have  $p_{k-1,i}^{\ell} \triangleright \mathbf{focus}(h_{k-1,i}, v) \iff p_{k,i}^{\ell} \triangleright \mathbf{focus}(h_{k,i}, v) \vee p_{k,i}^{\ell+1} \triangleright \mathbf{focus}(h_{k,i}, v)$ , and since the other patterns of  $\ell$  are unchanged,

$$(\forall i, p_{k-1,i}^{\ell} \triangleright \mathbf{focus}(h_{k-1,i}, v)) \iff (\forall i, p_{k,i}^{\ell} \triangleright \mathbf{focus}(h_{k,i}, v)) \vee (\forall i, p_{k,i}^{\ell+1} \triangleright \mathbf{focus}(h_{k,i}, v))$$

As the other rows are unchanged, the property holds for  $\mathcal{P}_k$ .

**Switch case**  $\mathcal{P}_{k-1} \rightarrow_{\text{switch},i,t} \mathcal{P}_k$  Let  $\mathcal{H}, f = \text{EXPAND}^\circ(\tau_i, t, h_{k,i})$ . Since the patterns outside of column  $i$  are unchanged, it is sufficient to prove that

$$p_{k-1,i}^{\ell} \triangleright \mathbf{focus}(h_{k-1,i}, v) \iff f(p_{k-1,i}^{\ell}) = \left(\overline{p_{k,i'}^{\ell'}}\right) \wedge \forall i' \in \mathcal{H}, p_{k,i'}^{\ell'} \triangleright \mathbf{focus}(h_{k,i'}, v)$$

We have  $p_{k-1,i}^{\ell} \triangleright \mathbf{focus}(h_{k-1,i}, v)$  iff.

1. their head constructors (if any) match and
2.  $\forall h_{k-1,i}[h] \in \mathcal{H}, \mathbf{focus}(h, p_{k-1,i}^{\ell}) \triangleright \mathbf{focus}(h_{k-1,i}[h], v)$ .

From the definitions of  $\text{EXPAND}^\circ$  and the sequence, we have

1. If  $f(p_{k-1,i}^{\ell}) = \emptyset$ , then:

$$\begin{aligned} p_{k-1,i}^{\ell} = c \wedge t \neq c \vee p_{k-1,i}^{\ell} = K(\dots) \wedge t \neq K &\iff \\ p_{k-1,i}^{\ell} = c \wedge \mathbf{focus}(h_{k-1,i}, v) = c' \wedge c' \neq c \vee p_{k-1,i}^{\ell} = K(\dots) \wedge \mathbf{focus}(h_{k-1,i}, v) = K'(\dots) \wedge K' \neq K & \end{aligned}$$

2. If  $f(p_{k-1,i}^{\ell}) \neq \emptyset$ , let  $\ell'$  be the row of  $\mathcal{P}_k$  yielded by  $f(p_{k-1,i}^{\ell})$ .

Then  $\forall h_{k,i'} = h_{k-1,i}[h] \in \mathcal{H}, p_{k,i'}^{\ell'} = \mathbf{focus}(h, p_{k-1,i}^{\ell})$  and thus:

$$p_{k,i'}^{\ell'} \triangleright \mathbf{focus}(h_{k,i'}, v) \iff \mathbf{focus}(h, p_{k-1,i}^{\ell}) \triangleright \mathbf{focus}(h_{k-1,i}[h], v)$$

And thus the property holds. □

### C.3.3 Bindings compatibility

For any pattern matrix of the sequence  $\mathcal{P}_k$  and any row  $\ell$  of this matrix, we define the set of contexts that designate variable subterms of the patterns of this row:

$$\begin{aligned} \mathbf{Var}_k^\ell &= \{h \mid \exists h_{k,i}, h'. h = h_{k,i}[h'] \wedge V(\mathbf{focus}(h', p_{k,i}^\ell)) \neq \emptyset\} \\ \text{where } V(p) &= \begin{cases} \{x\} & \text{if } p = x \text{ is a variable pattern} \\ V(q^1) & \text{if } p = q^1|q^2 \text{ and } V(q^1) = V(q^2) \\ \emptyset & \text{otherwise} \end{cases} \end{aligned}$$

**Lemma 5.** Let  $j$  an initial case. If  $p_0^j \triangleright v \rightarrow \sigma$ , then

$$\sigma = \left\{ x \mapsto \mathbf{focus}(h, v) \mid V(\mathbf{focus}(h, p_0^j)) = \{x\} \right\}$$

*Proof.* Trivial by induction on  $\triangleright$ .

**Lemma 6.** Let  $\mathcal{P}_k$  a pattern matrix in the sequence and  $\ell$  a row of  $\mathcal{P}_k$ . Let  $j = j_k^\ell$ . If  $p_{k,i}^\ell \neq -$ , then  $V(p_{k,i}^\ell) = V(\mathbf{focus}(h_{k,i}, p_0^j))$ .

*Proof.* By induction on  $k$

- if  $k = 0$ , we have  $V(p_0^j) = V(\mathbf{focus}(\square, p_0^j))$ .
- if  $k > 0$ , suppose the lemma holds for  $\mathcal{P}_{k-1}$ .
  - If  $\mathcal{P}_{k-1} \rightarrow_{\text{var}, i, \ell} \mathcal{P}_k$ , then  $p_{k,i}^\ell = \_$  and since the other patterns are unchanged, for every  $i', \ell'$  such that  $p_{k,i'}^{\ell'} \neq \_$ , we have:

$$V(p_{k,i'}^{\ell'}) = V(p_{k-1,i'}^{\ell'}) = V(\mathbf{focus}(h_{k,i}, p_0^{j_{k'}^{\ell'}}))$$

- If  $\mathcal{P}_{k-1} \rightarrow_{\text{or}} \mathcal{P}_k$  or  $\mathcal{P}_{k-1} \rightarrow_{\text{switch}} \mathcal{P}_k$ , then  $V(p_i^\ell)$  is unchanged for every  $i$  and  $\ell$  and the lemma still holds.  $\square$

**Lemma 7.** Let  $\mathcal{P}_k$  a pattern matrix in the sequence and  $\ell$  a row of  $\mathcal{P}_k$ . Let  $j = j_k^\ell$ . We have

$$s_k^\ell = \left\{ x \mapsto \text{KNIT}_{\tau_\Delta}^\bullet(h) \mid h \in \mathbf{Var}_0^j \setminus \mathbf{Var}_k^\ell \wedge V(\mathbf{focus}(h, p_0^j)) = \{x\} \right\}$$

*Proof.* By induction on  $k$

- if  $k = 0$ ,  $\ell$  is itself an initial case  $j$  and  $s_0^j = \emptyset = \mathbf{Var}_0^j \setminus \mathbf{Var}_0^j$ .
- if  $k > 0$ , suppose the property holds for every row of  $\mathcal{P}_{k-1}$ . If  $\mathcal{P}_{k-1} \rightarrow_{\text{or}} \mathcal{P}_k$  or  $\mathcal{P}_{k-1} \rightarrow_{\text{switch}} \mathcal{P}_k$ , the property still holds for every row as  $j^\ell$ ,  $s^\ell$  and  $\mathbf{Var}^\ell$  are unchanged.
  - If  $\mathcal{P}_{k-1} \rightarrow_{\text{var}, \ell, i} \mathcal{P}_k$ , let  $x = p_{k-1,i}^\ell$ . We have  $h_{k-1,i} \in \mathbf{Var}_{k-1}^\ell$ ,  $\mathbf{Var}_k^\ell = \mathbf{Var}_{k-1}^\ell \setminus \{h_{k-1,i}\}$  and  $s_k^\ell = s_{k-1}^\ell \cup \{x \mapsto \text{KNIT}_{\tau_\Delta}^\bullet(h_{k-1,i})\}$ . From Lemma 6, we have  $V(\mathbf{focus}(h_{k-1,i}, p_0^j)) = \{x\}$  and thus the property holds for the row  $\ell$ . It also holds for other rows, which are unchanged.  $\square$

**Proposition 2.** Let  $j$  an initial case. If  $p_0^j \triangleright v \rightarrow \sigma$  and the last pattern matrix of the sequence is  $\mathcal{P}_k$  with  $\text{COMPILE}_{\tau_\Delta}^\circ(\mathcal{P}_k) = \text{success}(j, s_k^0)$ , then we have

$$x \in s_k^0 \iff x \in \sigma \text{ and } \forall x \in \sigma, s_k^0(x)[r] \hookrightarrow \text{REPR}^\bullet(\sigma(x))$$

*Proof.* Since the sequence ends with  $\mathcal{P}_k$  (wildcard case), we have  $\mathbf{Var}_k^0 = \emptyset$ . The first part of the proposition then follows from Lemmas 5 and 7:

$$x \in s_k^0 \iff \exists h, V(\mathbf{focus}(h, p_0^j)) = \{x\} \iff (x \mapsto \mathbf{focus}(h, v)) \in \sigma$$

Let  $x$  be a variable bound in  $s_k^0$ . We have  $s_k^0(x) = \text{KNIT}_{\tau_\Delta}^\bullet(h)$  and  $\sigma(x) = \mathbf{focus}(h, v)$  for some context  $h$ .  $\text{KNIT}^\bullet$  guarantees that  $\text{KNIT}_{\tau_\Delta}^\bullet(h)[r] \hookrightarrow \text{REPR}_{\tau_\Delta}^\bullet(\mathbf{focus}(h, v))$ , thereby proving the property.  $\square$

### C.3.4 Termination

We define, for every matrix in the sequence  $\mathcal{P}_k$  and every row  $\ell$  of  $\mathcal{P}_k$ , the set of contexts such that the corresponding pattern in row  $\ell$  is not a wildcard:  $\mathbf{Pos}_k^\ell = \{h_{k,i} \mid p_{k,i}^\ell \neq \_ \}$ .

We also define a function similar to  $V$ :  $O(p) = \begin{cases} O(q^1) + O(q^2) + 1 & \text{if } p = q^1 | q^2 \\ 0 & \text{otherwise} \end{cases}$ .

Let  $\mathcal{O}_k = \sum \{O(p) \mid \exists h_{k,i}, h, \ell. p = \mathbf{focus}(h, p_{k,i}^\ell)\}$ .

**Proposition 3.** *Termination* There exists a  $k_0$  such that  $\mathcal{P}_{k_0}$  is the last pattern matrix of the sequence, i.e.,  $\text{COMPILE}_{\tau_\Delta}^\circ(\mathcal{P}_{k_0}) = \text{unreachable}$  (no pattern case) or  $\text{COMPILE}_{\tau_\Delta}^\circ(\mathcal{P}_{k_0}) = \text{success}(j, s)$  (wildcard case).

*Proof.* We show that there is a finite number of each kind of transition between pattern matrices of the sequence.

**Variable case** We trivially show that if  $\mathcal{P}_{k-1} \rightarrow_{\text{var}} \mathcal{P}_k$ , then  $\sum_{0 \leq \ell < m_k} |\mathbf{Var}_k^\ell| < \sum_{0 \leq \ell < m_{k-1}} |\mathbf{Var}_{k-1}^\ell|$ . Therefore the sequence contains at most  $\sum_{0 \leq j < m} |\mathbf{Var}_0^j| \rightarrow_{\text{var}}$ -steps.

**Or case** We trivially show that if  $\mathcal{P}_{k-1} \rightarrow_{\text{or}} \mathcal{P}_k$ , then  $\mathcal{O}_k < \mathcal{O}_{k-1}$ . Therefore the sequence contains at most  $\mathcal{O}_0 \rightarrow_{\text{or}}$ -steps.

**Switch case** We trivially show that if  $\mathcal{P}_{k-1} \rightarrow_{\text{switch}} \mathcal{P}_k$ , then  $\mathbf{Pos}_k^1 \subsetneq \mathbf{Pos}_{k-1}^1$  or  $m_k < m_{k-1}$  and therefore  $\sum_{0 \leq \ell < m_k} |\mathbf{Pos}_k^\ell| < \sum_{0 \leq \ell < m_{k-1}} |\mathbf{Pos}_{k-1}^\ell|$ . The sequence thus contains at most  $\sum_{0 \leq j < m} |\mathbf{Pos}_0^j|$ .

□

### C.3.5 Conclusion

Let  $\mathcal{P}_k$  be the last pattern matrix of the sequence.

- If  $\mathcal{P}_k = \emptyset$ , then  $\text{COMPILE}_{\tau_\Delta}^\circ(\mathcal{P}_k) = \text{unreachable}$  and from Proposition 1, we have  $p_0^j \not\triangleright v$  for every case  $j$ , that is  $p_0^0 \mid \dots \mid p_0^{m-1} \not\triangleright v$ .

- Otherwise, let  $j, s$  such that  $\mathcal{P}_k = \left( \begin{array}{ccc|c} h_{k,0} & \cdots & h_{k,n_k-1} & \\ \hline - & \cdots & - & j, s \\ \vdots & \ddots & \vdots & \vdots \\ p_{k,0}^{m_k-1} & \cdots & p_{k,n_k-1}^{m_k-1} & j_k^{m_k-1}, s_k^{m_k-1} \end{array} \right)$ .

From Proposition 1, we have  $p_0^j \triangleright v$  and for all  $j' < j$ ,  $p_0^{j'} \not\triangleright v$  (since  $\mathcal{P}^{j'} = \emptyset$ ), that is  $p_0^0 \mid \dots \mid p_0^{m-1} \triangleright v \rightarrow j$ .

Let  $\sigma$  such that  $p_0^0 \mid \dots \mid p_0^{m-1} \triangleright v \rightarrow j, \sigma$ . From Proposition 2, we have

$$s[r] \hookrightarrow \{x \mapsto \text{REPR}^\bullet(y)\}_{(x \mapsto y) \in \sigma}$$

which concludes.

□

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Algebraic Data Types and their memory representation</b>	<b>4</b>
<b>3</b>	<b>Pattern language</b>	<b>6</b>
3.1	Semantics . . . . .	6
3.2	Typing . . . . .	7
<b>4</b>	<b>Target Language and Memory Representations</b>	<b>8</b>
4.1	Memory values . . . . .	8
4.2	Computing on memory values . . . . .	8
4.2.1	Decision trees . . . . .	8
4.2.2	Tags and Contexts . . . . .	10
4.3	Memory representation . . . . .	10
4.4	A first example: the boxed representation . . . . .	11
4.4.1	Memory values . . . . .	11
4.4.2	Representation (Knit&Frog) . . . . .	12
<b>5</b>	<b>Representation-dependent compilation</b>	<b>13</b>
5.1	Algorithm . . . . .	14
<b>6</b>	<b>End-to-end examples</b>	<b>17</b>
6.1	Compilation procedure . . . . .	17
6.2	The OCaml representation . . . . .	19
6.2.1	Memory values . . . . .	19
6.2.2	Representation (Knit&Frog) . . . . .	19
6.2.3	Compilation . . . . .	20
6.3	The Packed Representation . . . . .	20
6.3.1	Memory values . . . . .	20
6.3.2	Representation (Knit&Frog) . . . . .	21
<b>7</b>	<b>Correctness</b>	<b>21</b>
7.1	Typing preservation for pattern matching . . . . .	22
7.2	Soundness . . . . .	22
<b>8</b>	<b>Related Works</b>	<b>22</b>
8.1	Memory representation and Algebraic Data Types . . . . .	22
8.2	Pattern languages . . . . .	23
8.3	Pattern Matching Optimization . . . . .	23
<b>9</b>	<b>Conclusion</b>	<b>23</b>
<b>A</b>	<b>Typing of values</b>	<b>29</b>
<b>B</b>	<b>Contexts and focusing</b>	<b>29</b>
<b>C</b>	<b>Proof of correctness</b>	<b>29</b>
C.1	Coherence of focusing . . . . .	29
C.2	Typing preservation . . . . .	29
C.3	Soundness . . . . .	31
C.3.1	Compilation steps . . . . .	31
C.3.2	Case matching equivalence . . . . .	31
C.3.3	Bindings compatibility . . . . .	32
C.3.4	Termination . . . . .	33

---

C.3.5 Conclusion . . . . .	34
----------------------------	----



**RESEARCH CENTRE  
GRENOBLE – RHÔNE-ALPES**

Inovallée  
655 avenue de l'Europe Montbonnot  
38334 Saint Ismier Cedex

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399