



HAL
open science

Compositional Flexible Memory Representations for Algebraic Data Types

Thaïs Baudon, Gabriel Radanne, Laure Gonnord

► **To cite this version:**

Thaïs Baudon, Gabriel Radanne, Laure Gonnord. Compositional Flexible Memory Representations for Algebraic Data Types. 9495, Inria. 2022. hal-03940742

HAL Id: hal-03940742

<https://hal.science/hal-03940742v1>

Submitted on 21 Aug 2024

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.



Distributed under a Creative Commons Attribution 4.0 International License

Inria

Compositional Flexible Memory Representations for Algebraic Data Types

Thaïs BAUDON, Laure GONNORD, Gabriel RADANNE

**RESEARCH
REPORT**

N° 9495

Decembre 2022

Project-Teams CASH

ISRN INRIA/RR--9495--FR+ENG

ISSN 0249-6399



Compositional Flexible Memory Representations for Algebraic Data Types

Thaïs BAUDON*, Laure GONNORD†, Gabriel RADANNE‡

Project-Teams CASH

Research Report n° 9495 — Decembre 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 *memory representation* which must accommodate garbage collection and polymorphism.

ADTs now appear in languages more liberal in their memory representation. Notably, Rust is now introducing more and more memory optimisations. As memory representation and compilation are interdependent, it raises the question of pattern matching compilation for highly customised layouts.

We propose to ease the experimentation of new, custom layouts in the early stages of compiler development by providing specification tools and a complete synthesis chain to generate pattern matching compilation procedures.

In this article, we present a novel way to specify *compositional* memory layouts, for which we automatically synthesise elementary builders and accessors, yielding a correct representation-specific compilation algorithm. This approach is implemented in a prototype tool **ribbit**.

Key-words: Algebraic Data Types, Memory Layouts, 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

RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES

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

Compositional Flexible Memory Representations for Algebraic Data Types

Résumé : Initialement présents dans les langages fonctionnels comme Ocaml et Haskell, les types algébriques sont maintenant de plus en plus présents dans les langages *mainstream* comme Rust. Ils permettent une abstraction commode des données et une façon élégante de décrire des fonctions via le filtrage de motifs (*pattern-matching*).

Ce rapport de recherche présente un framework facilitant la création et l'experimentation de nouvelle représentation mémoire en fournissant un outil de specification pour les représentations, et une chaine de synthèse completier pour générer les procedures de compilation attenante. La correction de l'algorithme est prouvée sous une hypothèse de *validité* de la représentation. L'algorithme est implémenté dans l'outil `ribbit`.

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

Contents

1	Introduction	4
2	Memory Representation of ADTs	5
2.1	Red-Black Trees in OCaml and Rust	5
2.2	Pattern Matching Optimising Compilation	6
2.3	Defining a Memory Representation	6
2.3.1	Composability	7
2.3.2	Distinguishability	7
2.3.3	Synthesis	7
3	Patterns and Pattern Matching	8
3.1	Patterns, Syntax and Semantics	8
3.2	Templated Compilation of Pattern Matching	9
4	Input Specification for Memory Layouts	12
4.1	Notations: Source and Memory Operations	12
4.2	Input Ingredients for Compositional Layouts	13
4.3	Example: OCaml Layout Specification	13
5	Synthesis of Memory Representations	14
5.1	Structural Interpolation	15
5.1.1	PREREPR_t^\bullet and REPRTY_t^\bullet	16
5.1.2	Distinguishability Constraints	16
5.1.3	Composability Constraints	16
5.2	Representation Completion	17
5.3	Synthesis of Pattern Matching Primitives	18
6	Extensions	19
7	Related Works	20
7.1	Memory Representation and ADTs	20
7.2	Pattern Matching Compilation	20
8	Conclusion	21
A	Source patterns	24
B	Memory patterns	24
C	Decision tree generation algorithm for $\text{INSPECTANDSPLIT}^\bullet$	25
C.1	Memory context relations	25
C.2	Decision tree generation algorithm	25
D	Full example: OCaml representation	27
D.1	Input specification	27
D.2	Primitive type t_{int}	27
D.3	Sum types t_{ABC} and t	27
D.4	Structural interpolation	28
D.5	Representation completion	30
D.6	Pattern matching primitives	31
D.6.1	Composability primitives: UNWRAP^\bullet	31
D.6.2	Distinguishability primitives: $\text{INSPECTANDSPLIT}^\bullet$	31
D.7	Pattern matching compilation	32

1 Introduction

Algebraic Data Types (ADTs) are an essential tool to model data. They allow to group together information in a consistent way through the use of records, also called product types, and to organise options through the use of variants, also called sum types. A proper ADT support enables to:

- 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.
- Optimise 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 for language designers wishing 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 optimisations 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 optimise the representation of terms. The simplest example is the `Option` type, whose values are either `Some value` or `None`. An easy way to represent such type is to box the value in the `Some` constructor below a pointer. However, if that value is an integer ranging from 0 to 10, we can represent it unboxed and use 11 for `None`. This optimisation is regularly done by programmers manually in for instance C++, at the price of error-prone manipulations. More complex optimisations on nested and rich data types are even more error-prone.

While these transformation are not generally easy, compilers have all the information at hand to perform them automatically on the whole program. This trove of potential optimisations has been brushed upon in recent versions of Rust, but remains largely unexplored. One reason for this lack of exploration is that ensuring the correctness of such optimising transformation is difficult. Indeed, the choice of memory representation of values in a language generally has far reaching consequences in the whole compiler, the language runtime but also it’s foreign function interface and garbage collector. Compiler designers are thus legitimately careful about changes in the memory representation, where mistakes could lead to complex compilation bugs.

Unfortunately, there is no ultimate memory representation: each language has very specific demands concerning its memory representation. Some languages requires great uniformity, for instance to allow for a tracing GC, other aim for the most compact representation. This prevent the development of a single highly efficient and formally proven memory representation that could be used in many languages.

We propose to address these difficulties from the perspective of a compiler designer that desires to explore various “memory optimisations” for monomorphic types and their performance regarding the pattern-matching procedure:

- We introduce a notion of *compositional* memory representation for Algebraic Data Types and define a novel specification for the compiler developer to provide such a representation.
- From this specification, and for each type of the input program, we automatically synthesise a function from values to their memory representation, along with a state-of-the-art optimising pattern matching compilation procedure tailored to the specified memory representation.
- Our synthesis procedure automatically checks that the provided specification is correct, and enforces the validity of the generated representation and the associated pattern matching compilation algorithm.
- Our framework is implemented in a tool `ribbit`.

```

//Colors for the Red-Black Trees
enum Color { Red, Black }

//Red-Black Trees of content T
enum RBT<T> {
  Node (Color, T, &RBT<T>, &RBT<T>),
  Empty,
}

pub fn cardinal<T>(v : RBT<T>) -> u32 {
  match v {
    Tree::Node(_, _, l, r) => 1+cardinal(*l)+cardinal(*r),
    Tree::Empty => 0
  }}

```

Figure 1: Red-Black trees and cardinal operation in Rust

```

fn cardinal(_1:RBT) = switch(_1){
  | 0 -> 0
  | 1 -> let _5 = cardinal*((_1 as Node).2))
        let _6 = cardinal*((_1 as Node).3))
        1 + _6 + _5
}

```

(a) Simplified Rust “MIR” output

```

cardinal (param/290) =
  (if (== param/290 1)
    0
    (+ (+ 1 (apply* cardinal (field2 param/290))
      (apply* cardinal (field3 param/290))))
  )

```

(b) Simplified OCaml “Lambda” output

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

2 Memory Representation of ADTs

To explore in more detail the memory representation of algebraic data types, we now investigate two languages which implements 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 favours performance and absolute control over low-level details. These differences result in drastically different choices in memory representation. We now focus on a concrete example of algebraic datatype which requires both expressivity and performance: Red Black Trees.

2.1 Red-Black Trees in OCaml and Rust

A Rust version of *Red-Black Trees* is depicted in Fig. 1. We first define `Color`, which is either the constant `Red` or `Black`. The type definition of `RBT<T>` expresses trees as recursive data structures of content `T`. It has two cases: `Empty`, and `Node` which contains a colour, a value, and two pointers (denoted `&(...)`) to its left and right subtrees. For instance, `Node(Red, 1515, &Node(Black, 0, &Empty,&Empty), &Empty)` is a tree of type `RBT<u64>`.

The `cardinal` function takes a tree and returns its cardinal 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, `Empty` yields a cardinal of 0 and `Node(_,_,l,r)` yields a cardinal of $1+\text{cardinal}(l)+\text{cardinal}(r)$.

We now look at how Rust and OCaml compile the `cardinal` 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

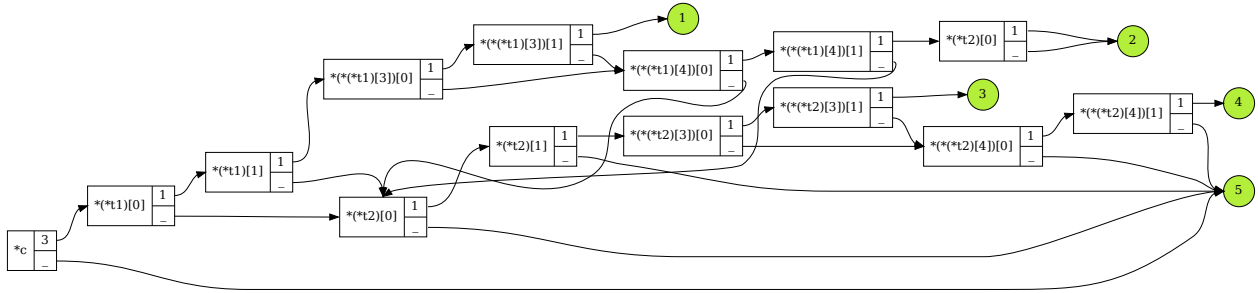


Figure 3: Output of our `ribbit` compiler for the balancing code of Fig. 1 in the OCaml memory representation. This output is a decision tree whose root is on the left. Decision (square) nodes compute values from their input (low-level representation), then branch depending on this value. Round nodes denote the final returned value of the procedure (the numbering of the matched branch).

representation. It uses the `if` primitive to discriminate between cases. Rust’s MIR¹ uses slightly lower level operations such as dereferencing and field accesses, and a `switch` construct on memory words.

In both cases, `Red`, `Black` and `Empty` are represented as unboxed integers. 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. In the `Node` case, the block contains a tag and four fields. This totals 6 memory words. In Rust, `Node` is directly a struct containing only four words. Indeed, Rust recognises that the tag (which differentiate `Node` and `Empty`) and the colour can be both packed in a single 64bits field. The resulting type uses 4 memory words.

2.2 Pattern Matching Optimising Compilation

Red-Black trees famously rely on a fairly complex balancing step, which redistributes the colours depending on the internal invariant of the data structure. Thanks to nested patterns, this step can be expressed very compactly using the following matching:

```
match c, v, t1, t2 {
  Black, z, &Node(Red, y, &Node(Red, x, a, b), c), d
| Black, z, &Node(Red, x, a, &Node(Red, y, b, c)), d
| Black, x, a, &Node(Red, z, &Node(Red, y, b, c), d)
| Black, x, a, &Node(Red, y, b, &Node(Red, z, c, d))
=> Node(Red, y, &Node(Black,x,a,b), &Node(Black,z,c,d)),
  a, b, c, d => Node(a, b, c, d),
}
```

This pattern matching inspects four arguments at the same time: the current colour `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 optimised decision trees Kosarev et al. (2020); Maranget (2008); Sestoft (1996). The resulting code is highly non-trivial, as can be seen in Fig. 3.

2.3 Defining a Memory Representation

So far, we have kept a prudish veil on how one actually defines the memory representation of values in a language like Rust or OCaml. This is largely because it is composed of an uncountable collection of decision spread around the whole design of the compiler and its runtime: it affects the compilation of expressions, creation of values, field accesses, foreign function interfaces, runtime and garbage collection code and, what concerns us most in this article, pattern matching compilation. All these elements must naturally be kept synchronised, at the risk of causing delicate compiler bugs. Even in a language like Rust, which favour optimised memory layouts, incremental improvements to the representation of values has progressed slowly due to the far-reaching nature of such changes to the compiler.

¹the Middle Intermediate Representation (MIR) is a control flow graph. We reconstruct some of its logic here for ease of reading.

In this article, we present a framework to easily define the memory representation of a language. Our framework automatically checks, using constraint solving, that the proposed representation is well-formed – for instance, that constructors can effectively be differentiated – and synthesise two functions: one to compile pattern matches and one to turn values into their memory representation.

Our framework is based on two main insights.

1. A reasonable memory layout of values should be **composable**: for a value v , its subvalues can be found in its representation. For instance `Node(Red, 1515, &Node(Black, 0, &Empty,&Empty), &Empty)` contains the subvalue `Node(Black, 0, &Empty,&Empty)` in its representation.
2. To yield itself to pattern matching, the memory representation of values should be **distinguishable**: one should be able to determine if an RBT value is a `Node` or an `Empty`.

To specify such a memory representation, called *compositional*, one should provide two operations for each type: 1. where to place subterms in the memory representation, 2. how to distinguish values. For simplicity, we only show the definition for the `RBT<u64>` type. A full compiler would provide such information for all types.

2.3.1 Composability

For each type, we require the position of its subterm. Concretely, this is a mapping from source values with holes to memory values with holes:

$$\begin{aligned}
 \text{Node}(c : \text{Color}, _ , _) &\rightarrow \{\{\mathcal{L}_{64}, _ , _ , _ \}\} \text{ with } c'[0, 1] = 2 * c + 1 \\
 \text{Node}(_ , i : \text{u64}, _ , _) &\rightarrow \{\{_ , \mathcal{I}_{64}, _ , _ \}\} \\
 \text{Node}(_ , _ , l, _) &\rightarrow \{\{_ , _ , \mathcal{L}_*, _ \}\} \\
 \text{Node}(_ , _ , _ , r) &\rightarrow \{\{_ , _ , _ , \mathcal{R}_*\}\} \\
 \&t &\rightarrow \text{Ptr}_{64,3}\mathcal{L}_*
 \end{aligned}$$

There are many things to note here. $\{\{\ \}\}$ indicates a structure: a sequence of memory word of various sizes. $_*$ is an optionally sized position in the representation. Some fields are directly sized: the integer content i for instance is explicitly sized as 64 bit. The sizes of l and r , the two subtree pointers, are left unspecified and will be inferred. Numerical representation can be tweaked as we move to the memory representation: for instance, the colour c is stored in the first two bits of the header word, appropriately shifted. We also indicate that the pointer is 64-bits wide, and that 3-bits are always 0 for alignment.

To enforce composability, we enforce that the “source” patterns on the left hand side should be *shallow* – without any subpatterns. Note that there are no mapping for `Empty`, since it doesn’t have subvalues.

2.3.2 Distinguishability

To distinguish values, we require a collection of *discriminants* that express where in memory to look, and which value would be associated to each constructor. We only need one discriminant for `RBT`:

$$\{\{\mathcal{I}_{64}, \dots\}\} \text{ with } \text{tag} = h[0, 0] \Rightarrow \begin{cases} \text{Empty} \mapsto \text{tag} = 0 \\ \text{Node} \mapsto \text{tag} = 1 \end{cases}$$

This reads as follow: “Let tag be the lowest bit of the first field of the given RBT value, if the constructor of the value is `Empty` it has $\text{tag} = 0$, if it is `Node` is has $\text{tag} = 1$ ”.

2.3.3 Synthesis

These two information are in fact completely sufficient to fully determine the memory representation of values. Our synthesis proceeds as follow:

1. We first derives a “general shape” of values. Here, this indicate that `RBT` memory values are of the shape $\{\{\text{Word}_{64}\}\}$ in the `Empty` case – a struct with a single 64bits integer – and $\{\{\text{Word}_{64}, \text{Word}_{64}, \text{Ptr}_{64,3}\text{RBT}, \text{Ptr}_{64,3}\text{RBT}\}\}$ in the `Node` case – a struct with four 64bits words.

2. Based on this general shape, the missing information is inferred using constraint solving. For instance, it will determine that the word in the *Empty* case will always be 0, to distinguish with the *Node* case. This step also checks the correctness of the representation.
3. Thanks to this last step, the memory representation is fully defined. The framework synthesise a REPR^\bullet functions which returns the memory value of any given value. For instance $\text{REPR}_{RBT\langle u64 \rangle}^\bullet(\text{Node}(\text{Red}, 1515, \text{Empty}, \text{Empty}))$
 $\{\{2, 1515, \text{Ptr}_{64,3} \{0\}, \text{Ptr}_{64,3} \{0\}\}\}$
4. Finally, our framework synthesise a compilation function from pattern to decision trees.

Our synthesis framework has been implemented in a tool called `ribbit`. It provides the full chain from specification of the memory representation (for now, as OCaml code), to constraint solving (using the SMT solver `Z3`) to pattern matching compilation. We showcase the compilation of the balancing function provided above in Fig. 3.

In the rest of this article, we formally define a language with pattern matching (Section 3), propose an input specification for *compositional representations* (Section 4) and details our synthesis procedure (Section 5). We finally sketch some extensions to our framework (Section 6).

3 Patterns and Pattern Matching

In this section, we present a language of patterns as well as a *generic* algorithm to compile pattern matching to decision trees.

3.1 Patterns, Syntax and Semantics

The PAT language, described in Fig. 4, 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 constructor patterns along with variables (x), wildcards ($-$) and “or”-patterns ($p_1 \mid p_2$). Type expressions, denoted τ , are composed of sums, written $\sum_{0 \leq i < n} K_i(t_{i,0}, \dots, t_{i,n_i-1})$ or $K_0(\tau_{0,0}, \dots, \tau_{0,n_0-1}) + \dots + K_{n-1}(\tau_{n-1,0}, \dots, \tau_{n-1,n_{n-1}-1})$, and of primitive types (essentially fixed-width signed or unsigned integers). Type environments, denoted Γ , contains both variable and type names, both associated to type expressions. Values, denoted v , follow a subset of the pattern grammar. Finally, value environments, denoted σ , associate variables to their values.

For ease of presentation, we assume the following:

- Types are monomorphic: all code has been specialised.
- All types are named with unique identifiers, whose definitions are stored in the type environment Γ . This allows us to segment types and break recursions. For instance the definition $\tau = A + B(i32)$ yields the type naming environment $\Gamma = \{t \mapsto A + B(t'); t' \mapsto i32\}$. Recursive types result in a cycle in the type environment.
- Patterns are *exhaustive*: all possible cases are handled.

Dynamic Semantics The semantics of patterns take the form of a judgement $p \triangleright v \rightarrow \sigma$ meaning that pattern p matches value v and binds the variables bound in σ . We also define the matching judgement $\text{Match} \{p_1 \mid \dots \mid p_n\} \triangleright v \rightarrow i, \sigma$ which additionally returns the index of the matched branch. In a full-fledged language, it would then trigger the evaluation of the body of the branch in question. The semantic rules are standard and shown in Appendix A.

Static Semantics The typing judgement, denoted $\Gamma \vdash p : t \rightarrow \Gamma'$, indicates that pattern p is of type t (defined in Γ) and returns a new environment Γ' containing the new variables bound in p . It follows the standard rules for typing patterns, shown in Appendix A.

Example 1 (Running example: definition). As a running example for the rest of the paper, we consider the type t defined in the following type naming environment:

Patterns

$p ::= _$	(Wildcard)
$\quad x$	(Variable)
$\quad (p_1 \mid p_2)$	(Disjunction)
$\quad K(p_1, \dots, p_n)$	(Constructor pattern)

Matching

$m ::= \text{Match } \{p_1 \mid \dots \mid p_n\}$	(Matching)
---	------------

Types

$t ::= T$	(Primitive type)
$\quad \sum_n K_i(t_0, \dots, t_{n_i-1})$	(Sum type)
$\Gamma ::= \{x_i \mapsto \tau_i; t_j \mapsto \tau_j\}$	(Type env.)

Values

$v ::= z \in \mathbb{Z}$	(Integer constant)
$\quad K(v_0, \dots, v_{n-1})$	(Constructor value)
$\sigma ::= \{x_0 \mapsto v_0; \dots; x_{n-1} \mapsto v_{n-1}\}$	(Value env.)

Figure 4: PAT, our simplified language of patterns and types

$$\Gamma = \left\{ \begin{array}{l} t \mapsto \text{None} + \text{Some}(t_{\text{ABC}}); \\ t_{\text{ABC}} \mapsto A + B + C(t_{\text{int}}); \quad t_{\text{int}} \mapsto \text{u32} \end{array} \right\}$$

We define the following matching problem:

$$m_0 = \text{Match} \left\{ \begin{array}{l} | (p_0) \text{None} \mid \text{Some}(A) \\ | (p_1) \text{Some}(B) \\ | (p_2) \text{Some}(C(n)) \end{array} \right\}$$

- The *pattern* p_0 matches either value *None* or *Some*(A). We write $m_0 \triangleright \text{None} \rightarrow 0, \emptyset$ and $m_0 \triangleright \text{Some}(A) \rightarrow 0, \emptyset$.
- The third pattern p_2 matches any value of type t of the form *Some*($C(n)$), where n is a value of type t_{int} , and binds its value to the given name. For instance, $m_0 \triangleright \text{Some}(C(42)) \rightarrow 2, \{n \mapsto 42\}$

Additionally, these patterns are well-typed:

$$\Gamma \vdash p_0 : t \rightarrow \Gamma \qquad \Gamma \vdash p_2 : t \rightarrow \Gamma \cup \{n : \text{u32}\}$$

◆

3.2 Templated Compilation of Pattern Matching

We now quickly rephrase Maranget (2008)'s state-of-the-art pattern matching compilation algorithm to decision trees, illustrated in Example 3. Decision trees $\mathcal{T} \in \text{Trees}$, whose grammar is described in Fig. 5, formalise the description of *nested switch cases* that we illustrated in Fig. 3.

The function $\text{COMPILE}_t : \text{PAT} \rightarrow \text{Trees}$ takes as input a pattern matching problem for values of a given type t and outputs an equivalent decision tree for *memory representations* of such values. In the rest of this article, we assume that Γ is always provided implicitly.

Decision trees Decision trees are composed of nodes of the form $\text{switch}(e) \{ \mathcal{C} \}$, consisting of an expression e computing the memory value under scrutiny and of a list of cases \mathcal{C} . A decision tree leaf is of the form

Expressions

$e ::= \Delta$	(Main input)
$e \oplus z$ where $\oplus \in \{\wedge, \vee, +, \leq, \neq, \dots\}$	(Bitwise and arithmetic operations)
$* e$	(Pointer dereferencing)
$e.i$	(i -th block field access)
\dots	(Memory-dependent operations)

Switch cases

$\mathcal{C} ::= z \mapsto \mathcal{T} \cdot \mathcal{C}$	(Regular case)
$\top \mapsto \mathcal{T}$	(Default case)
\emptyset	(No default case)

Decision trees

$\mathcal{T} \in \text{Trees}$	
$\mathcal{T} ::= \text{switch}(e) \{ \mathcal{C} \}$	(Decision node)
$\text{success}(j, \{\overline{x_i \mapsto e_i}\})$	(Branch j with bindings)
unreachable	(Unreachable leaf)

Figure 5: Decision trees and their expressions

$\text{success}(j, \sigma)$, which successfully returns a branch index j and a set of binding expressions σ . Since we only consider exhaustive patterns, there are no failure cases.

At toplevel, decision trees take as input the *main input* of the whole matching, denoted Δ . The scrutinee of a switch is an expression manipulating memory values (thus an expression over Δ). These expressions are partially target-dependent but composed at least of dereferencing, field access and simple bitwise and arithmetic operations. The output bindings σ are also target-dependent.

The dynamic semantics of these decision trees are similar to those of pattern matching, thus we omit their description and provide a simple example.

Example 2 (Decision tree). The following decision tree matches an option type in the OCaml representation:

$$\text{switch}(\Delta \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \text{success}(0, \emptyset) \\ 0 \mapsto \text{success}(1, \{n \mapsto * \Delta.1\}) \end{array} \right\}$$

This first determines whether the least significant bit of Δ is 1 or 0. This determines whether a value is *None* (the integer 1) or a *Some* (a pointer). If *None*, we succeed with branch 0. If *Some*, we succeed with branch 1 and a *binding environment* mapping n to the value inside *Some*, as computed by $* \Delta.1$. \blacklozenge

Pattern matching compilation algorithm Our algorithm $\text{COMPILE}_{t_0}(\mathcal{P})$ compiles a pattern matching problem to an equivalent decision tree using the now classic approach of *pattern matrices*. Its arguments are:

- the type of the main input value (at toplevel) t_0 ;
- a *pattern matrix* \mathcal{P} in which each row encodes a case of the matching problem and each column a term to be scrutinised (allowing us, for instance, to split product patterns and choose which field to inspect first). At the end of each row, we also record its index and a binding environment. The column headers indicate which parts of the main input this column is matching against, as a subterm of the toplevel value.

COMPILE_t proceeds by recursively emitting switch nodes (depicted as branching nodes in our graphical representation) and reducing the pattern matrix until exhaustion. It has five cases, depending on the shape of \mathcal{P} :

empty If \mathcal{P} is empty, the algorithm fails, as exhaustive patterns should have a valid pattern for every case.
wildcard If the first row consists only of wildcards, the matching can only succeed as this branch accepts all values.

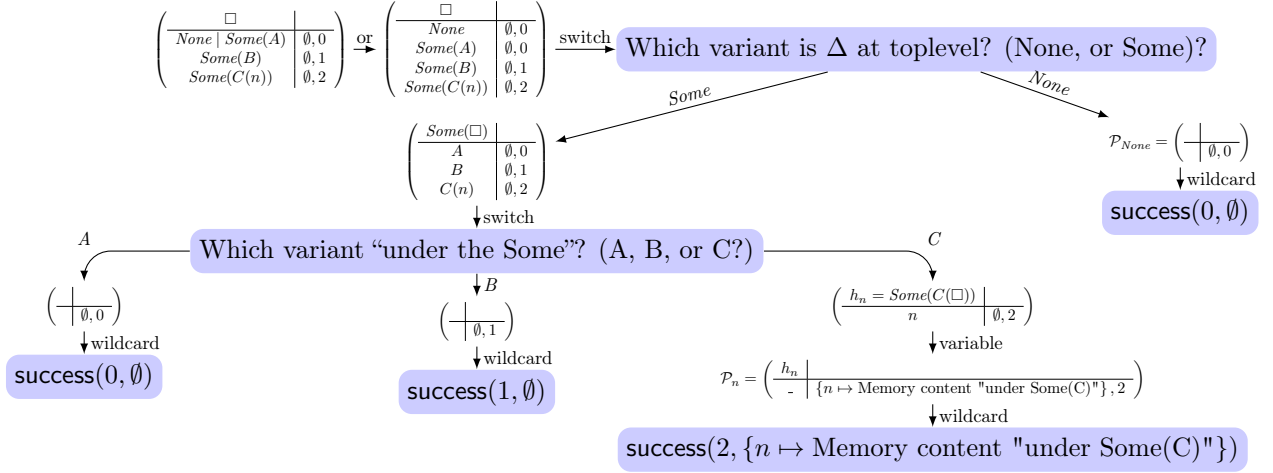


Figure 6: Execution of $\text{COMPILER}_t(m_0)$ on the running example. Blue coloured parts depict the final decision tree (that can be used to match any value of type t); other parts depict pattern matrices used during the compilation procedure.

variable If there exists a variable pattern p_i^ℓ in some row ℓ (and the previous rows contain no variable patterns), we introduce p_i^ℓ into this row’s binding environment.

or If there exists an “or”-pattern $p_i^\ell = p_1 \mid p_2$ in some row ℓ (and the previous rows contain no “or”-patterns), we split it into two new rows.

switch Otherwise, an actual switch node must be generated. The aim of this switch node is to inspect the head constructor of the value contained in one of the inputs, compare it to the patterns of its column in \mathcal{P} , then branch to its associated subtree that performs the remaining computations (on other columns and nested values).

Example 3 (Running example: $\text{COMPILER}_{t_\Delta}$ function). Fig. 6 depicts the execution of $\text{COMPILER}_t(m_0)$ on the running example (Example 1). The produced decision tree can then be used to “match” the memory representation of any value v of type t . It first applies the **or** rule, then the **switch** rule to decide whether the head constructor of v is *None* or *Some*. In the former case, it immediately stops, returning 0. In the latter case, it continues through the left branch of the tree while getting access to the subvalue “under *Some*”, for which it should decide whether it is *A*, *B* or a *C*. Finally, if this subvalue is a *C*, it should get access to the subvalue “under *C*”. Note that the detailed operation of how to discriminate between constructors such as *A*, *B* and *C* is still undefined at this point. This is indeed highly dependent on precise details of the chosen internal memory layout. \blacklozenge

Representation dependence While Maranget (2008)’s algorithm produces good decision trees, it was originally specific to a given memory representation (the OCaml one). In order to allow for custom representations, we must provide additional representation-specific functions to:

- extract the representation of a subvalue from the memory representation of its parent value (UNWRAP_t^\bullet);
- retrieve the head constructor of a source value from its memory representation ($\text{INSPECTANDSPLIT}_t^\bullet$).

Unfortunately, specifying these functions for each type and constructor is particularly complex. In particular, ensuring the correctness of such functions when written by hand in compiler code is delicate. Instead, we propose to *synthesise* such representation-specific functions for any custom memory representation, based on a higher-level description of the memory layout of values. In the following sections, we show how to automatically derive the *full compilation algorithm for any type*, including representation-specific functions, from a user-provided *high-level description* of the desired memory representation.

$$\begin{aligned}
\widehat{\tau} \in \widehat{Types} &::= - \mid \text{Word}_\ell \mid \text{Ptr}_{\ell,a}\widehat{\tau} \mid \left\{ \left\{ \widehat{\tau} \right\} \right\} \mid \sqcup \widehat{\tau} \\
\widehat{v} \in \widehat{Values} &::= - \mid \text{Word}_\ell(z) \mid \text{Ptr}_{\ell,a}\widehat{v} \mid \left\{ \left\{ \widehat{v} \right\} \right\} \mid \sqcup \widehat{v} \\
\text{ops} \in \text{WordOps} &::= \blacksquare \quad (\text{Identity}) \\
&\mid \text{ops}[i,j] \quad (\text{Bits } i \text{ to } j, \text{ inclusive}) \\
&\mid \text{ops} \oplus z \quad \text{where } \oplus \in \{\wedge, \vee, +, \leq, \neq, \dots\} \quad (\text{Bitwise and arithmetic operations}) \\
\widehat{h} \in \widehat{Contexts} &::= \text{ops}_* \quad (\text{Unsize memory hole}) \\
&\mid \text{ops}_\ell \quad (\text{Sized memory hole}) \\
&\mid \text{Ptr}_{\ell,a}\widehat{h} \quad (\text{Pointer}) \\
&\mid \left\{ \left\{ -, \dots, \widehat{h}_i, \dots, - \right\} \right\} \quad (i\text{-th field in a block})
\end{aligned}$$

Figure 7: Memory values, types and contexts

4 Input Specification for Memory Layouts

A layout, or *memory representation*, specifies how to encode values of any source type in memory, so as to be able to encode each typed value in a fully deterministic way. This section formally explains how to provide an input specification for our synthesis algorithm.

4.1 Notations: Source and Memory Operations

In this section, we define similar notions on the source language PAT and on memory contents. From now on, hats will denote memory-related notions: $h \in \text{Contexts}$ denotes a *source* context while $\widehat{h} \in \widehat{\text{Contexts}}$ denotes a *memory* context.

Source contexts, subterms of a type Source contexts are used to deconstruct values and types. A constructor context is a constructor pattern in which all subpatterns are wildcards, except one which is another context. All contexts end with a single *hole* \square . The grammar for source contexts (*Contexts*) is defined below:

$$\begin{aligned}
h &::= \square \quad (\text{Hole}) \\
&\mid K(_, \dots, _, h, _, \dots, _) \quad (\text{Constructor context})
\end{aligned}$$

We can now define the *subterms* of a sum type expression as a set of contexts where subterms are located:

$$\text{SubTerms}\left(\sum_n K_i(\overline{t_{i,j}^{n_i}})\right) = \left\{ K_i\left(-, \dots, \square_j, \dots, -\right) \right\}_{\substack{0 \leq i < n \\ 0 \leq j < n_i}}$$

We also write $\text{SubTerms}(t)$ instead of $\text{SubTerms}(\Gamma(t))$ when Γ is implicit from the context.

Example 4 (Running example: subterms). From the definition of types in Example 1, we obtain:

$$\text{SubTerms}(t) = \{\text{Some}(\square)\} \quad \text{SubTerms}(t_{\text{ABC}}) = \{C(\square)\}$$

Constant constructors such as *None* do not yield any subterm context. ◆

Memory Values and Types In order to work on the *internal representation* of values, we consider an abstraction of memory, similar to LLVM IR’s types, shown in Fig. 7. A *memory value* is always typed and is either a fixed-size word encoding an integer $z \in \mathbb{Z}$, written $\text{Word}_\ell(z)$ and typed Word_ℓ , a fixed-size pointer to another memory value \hat{v} of type $\hat{\tau}$ with a fixed number a of address alignment bits (filled with 0s), written $\text{Ptr}_{\ell,a}\hat{v}$ and typed $\text{Ptr}_{\ell,a}\hat{\tau}$, or a contiguous block containing a finite number of heterogenous memory values (fields), written $\{\{\hat{v}_0, \dots, \hat{v}_{n-1}\}\}$ and typed $\{\{\hat{\tau}_0, \dots, \hat{\tau}_{n-1}\}\}$. We also define a (finite) *disjoint union* of memory values (resp. types) to denote the possibility of several concrete memory values (resp. types) at a given location, along with a wildcard $_$ to denote an unspecified memory value (resp. type). The typing rules for the judgment $\vdash \hat{v} : \hat{\tau}$ are straightforward and omitted for the sake of brevity.

Memory Contexts A *memory context*, defined in Fig. 7, indicates a position within a memory value and an optional transformation of the word at this position. It consists of a coarse-grained hierarchy of *structural* elements – namely pointers and blocks – ending with an optionally sized “hole” $_ \ell$ and of finer-grained operations on the word designated by the hole. The hole of a memory context applied to a word (resp. pointer) value should match the size of this word (resp. the number of alignment bits). The hole of a memory context applied to a block must be unsized and contain the identity operation (i.e., \blacksquare_*).

Example 5 (Running example: memory contexts). Consider t_{ABC} from our running example (Example 1). Its memory representation in OCaml (which we will refine later in Section 4.3) is of type $\hat{\tau} = \text{Word}_{64} \sqcup \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{Word}_{64}\}\}$. The OCaml representation maps the concrete source values A , B and $C(42)$ to the following memory values:

$$A \mapsto \text{Word}_{64}(1) \quad B \mapsto \text{Word}_{64}(3) \quad C(42) \mapsto \text{Ptr}_{64,3} \{\{\text{Word}_{64}(2), \text{Word}_{64}(85)\}\}$$

The two constant constructors A and B are mapped to odd constant values. The C constructor is mapped to a 64-bit-wide, 3-bit-aligned pointer (thus an even number), pointing to a block containing its tag 2 and a value. The integer value 42 is located in the second field of this block. As we will see in more detail in Section 4.3, integers are encoded through a 1-bit left shift and an increment ($42 \times 2 + 1 = 85$). \blacklozenge

4.2 Input Ingredients for Compositional Layouts

Let t a source type such that $\Gamma(t) = \sum_n K_i(\overline{t_{i,j}}^{n_i})$. In order to synthetise a full (compositional) memory representation, the user should provide:

- $\text{POSITION}_t^\bullet : \text{SubTerms}(t) \rightarrow \widehat{\text{Contexts}}$
defines how the representation of each subvalue is composed into the representation of its parent value;
- $\text{DISCRS}_t^\bullet \subset \widehat{\text{Contexts}} \times (\{K_i\}_n \rightarrow \mathbb{Z})$
is a set of *discriminants*. A discriminant consists of a memory context indicating where to look, and a *partial* mapping from constructors to integer values that can be found at this position. The discriminants should be defined such that, for each distinct pair of constructors, at least one discriminant is defined for both constructors and maps them to distinct integer values, i.e.:

$$\forall 0 \leq i < j < n, \exists (\hat{h}, \text{split}) \in \text{DISCRS}_t^\bullet, \text{split}(K_i) \neq \text{split}(K_j)$$

4.3 Example: OCaml Layout Specification

We now demonstrate how to specify the OCaml memory representation of values Minsky and Madhavapeddy (2021) in our framework. All computation details can be found in Appendix D.

Informal specification OCaml memory values are either unboxed immediates or pointers to blocks, both one word wide. OCaml blocks are contiguous word-aligned structs. Every block starts with a *header* word containing various metadata including a *tag* indicating the constructor of the underlying value. For simplicity, we ignore the remaining metadata here. An important optimisation is that argument-less constructors are represented as unboxed integers, rather than empty blocks. The integer value is the index of the constructor. In OCaml, *None* has the same representation as the integer value 0 (i.e., $\text{Word}_{64}(1)$), without any boxing.

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 pointer. Finally, addresses are word-aligned, guaranteeing that their 3 last bits are always 0 (for 64-bit-wide words).

Formal OCaml layout specification Let t such that $\Gamma(t) = \sum_n K_i (\overline{t_{i,j}^{n_i}})$.

The POSITION^\bullet function gives the position of a subvalue in the “pointer-to-block” memory structure:

$$\text{POSITION}_t^\bullet \left(K_i \left(_, \dots, \square, \dots, _ \right) \right) = \text{Ptr}_{64,3} \left\{ \left\{ _, \dots, \blacksquare_{j+1}^*, \dots, _ \right\} \right\}$$

The three discriminants express how to distinguish between constructors:

- $D_{\geq 0}$ distinguishes between unit and non-unit constructors by inspecting the lowest bit, which is always 1 (resp. 0) for unit (resp. non-unit) constructors:

$$D_{\geq 0} = \left(\blacksquare[0,0]_{64}^* ; \left\{ K_i \mapsto \begin{cases} 1 & \text{if } n_i = 0 \\ 0 & \text{otherwise} \end{cases} \mid 0 \leq i < n \right\} \right)$$

- $D_{=0}$ distinguishes individual unit constructors:

$$D_{=0} = (\blacksquare_{64} ; \{ K_i \mapsto 2 \times i + 1 \mid n_i = 0 \})$$

In this case, the value must be an unboxed integer, which we can inspect right away.

- $D_{>0}$ distinguishes individual non-unit constructors:

$$D_{>0} = (\text{Ptr}_{64,3} \{ \blacksquare_{64}, \dots \} ; \{ K_i \mapsto i \mid n_i > 0 \})$$

In this case, the value must be a pointer to a block.

We then have $\text{DISCRS}_t^\bullet = \{ D_{\geq 0}, D_{=0}, D_{>0} \}$.

Example 6 (Running example: t_{ABC} ingredients). Consider t_{ABC} in our running example. As we saw in Example 5, the unit constructors A and B are distinguishable from C by inspecting the lowest bit of the representation (which acts as a pointer flag); we also define a discriminant to distinguish A from B and another that solely affects C (so as to obtain the expected memory structure).

$$\begin{aligned} \text{DISCRS}_{t_{ABC}}^\bullet = & \left\{ \left(\blacksquare[0,0]_{64}, \begin{cases} A \mapsto 1 & C \mapsto 0 \\ B \mapsto 1 \end{cases} \right), \right. \\ & (\blacksquare_{64}, \{ A \mapsto 1 \quad B \mapsto 3 \}), \\ & \left. (\text{Ptr}_{64,3} \{ \blacksquare_{64}, \dots \}, \{ C \mapsto 2 \}) \right\} \end{aligned}$$

We also define the composition of $\text{REPR}_{t_{\text{int}}}^\bullet(n)$

within $\text{REPR}_{t_{ABC}}^\bullet(C(n))$:

$$\text{POSITION}_{t_{ABC}}^\bullet(C(\square)) = \text{Ptr}_{64,3} \{ _, \blacksquare_{64}^*, \dots \}$$

5 Synthesis of Memory Representations

We now give our complete synthesis procedure, whose overall pipeline is shown in Fig. 8, which takes as input the specification functions of Section 4.2.

It proceeds by generating partial intermediate representations of a type t for each tag , which can be either a constructor or a wildcard \top . It then completes these intermediate representations using constraint solving. This is done in three steps:

1. *Structural interpolation* (Section 5.1) determines the memory type $\text{REPRTY}_t^\bullet(tag)$ of the representation of any value of type t whose constructor fits tag , as well as a *memory pattern* $\text{PREREPR}_t^\bullet(tag)$ which encodes all statically known representation elements. This step proceeds by collecting all necessary structural information from the input specification and computing an “intersection” of all collected constraints. If this step succeeds, its result is a memory pattern with some remaining unsolved word contents.

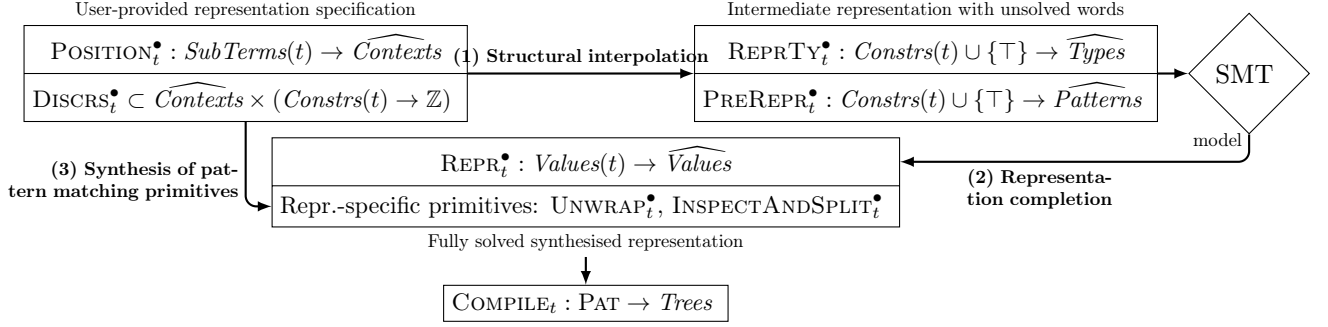


Figure 8: Representation synthesis pipeline

$$\begin{aligned} \phi &::= \top \mid (\text{ops} = x_h) \mid (\text{ops} = z) \mid \phi \wedge \phi \\ \widehat{p} \in \widehat{Patterns} &::= t@h \mid \text{Word}_\ell(\phi) \mid \text{Ptr}_{\ell,a}\widehat{p} \mid \left\{ \left\{ \widehat{p}^? \right\} \right\} \mid \bigsqcup \widehat{p} \end{aligned}$$

Figure 9: Memory patterns and word formulas

2. *Representation completion* (Section 5.2) fully determines the contents of each word. After ensuring that the intersection of all patterns is non-empty, there still remain under-specified parts. The goal of this step is to synthesise a *unique memory value*. We translate unsolved word descriptions to SMT formulas on fixed-size bitvectors. Solving these formulas yields a model that assigns a concrete value or a subterm symbol to each word, which finishes the complete synthesis of the REPR_t^\bullet function. If a formula is unsatisfiable, it means that we cannot obtain a suitable representation.
3. *Synthesis of pattern matching primitives* (Section 5.3) produces the primitives needed by our templated COMPILE_t algorithm. More precisely, for each source type t , we generate a function UNWRAP_t^\bullet that extracts the representation of a subvalue from the representation of its parent value and $\text{INSPECTANDSPLIT}_t^\bullet$, which determines the constructor of a source value from its representation.

Each of these three steps may fail, indicating that the input specification is *invalid*. The rest of the section follows the synthesis pipeline. Since the representation of a type is dependent on the representation of its subtypes, we process the types of an environment Γ in the reverse prefix order.

5.1 Structural Interpolation

$\text{PREREPR}_t^\bullet(\text{tag})$ is representation template in which structural elements are fully determined but word contents remain unsolved. $\text{REPRTY}_t^\bullet(\text{tag})$ gives the memory type for a given tag. To generate them, we use *memory patterns*, shown in Fig. 9. Memory patterns are memory value skeletons in which word contents are specified by a formula ϕ . A formula is either true (\top) or a conjunction of clauses of the form $(\text{ops}_1 = x_1) \wedge \dots \wedge (\text{ops}_m = x_m)$, where x_k are concrete integers or subterm symbols and ops_k are sequences of operations applied to the word contents \blacksquare . Memory patterns also feature an application case $t@h$, which stands for the (unsolved) memory representation of the subterm of type t at position h . The remaining pattern cases are similar to the previously defined memory values and types.

Memory patterns are equipped with *union* and *intersection* operations, defined in Appendix B, which allow to build them up incrementally. For instance, the following rule expresses that a pointer pattern is compatible with any word pattern whose width “fits” in its address alignment bits; the specified word contents must then be zero.

$$\frac{\ell' \leq a \quad (0)_{\ell'} \text{ satisfies } \phi}{\text{Ptr}_{\ell,a}\widehat{p} \cap \text{Word}_{\ell'}(\phi) = \text{Ptr}_{\ell,a}\widehat{p}}$$

Memory patterns also have an *instantiation* operation, noted $\widehat{h}[\widehat{p}]$ which builds a pattern from \widehat{h} , where the hole is replaced by \widehat{p} . This is also defined in Appendix B.

Example 7 (Memory patterns). $\text{Word}_{64}(\blacksquare[0,0] = 1)$ denotes any 64-bit-wide word whose lowest bit is 1. In our running example, the memory pattern $t_{\text{int}}@C(\square)$ matches memory values representing a value of type t_{int} . Ultimately, it will represent $\text{REPR}_{t_{\text{int}}}^\bullet(v)$ where v is the subvalue “under C ”. \blacklozenge

5.1.1 PREREPR_t^\bullet and REPRTY_t^\bullet

Let t a type and tag one of its constructors, or \top . $\text{PREREPR}_t^\bullet(tag)$ is the intersection of all constraints pertaining to this memory layout: composability constraints, denoted \widehat{p}_t^c and distinguishability constraints, denoted \widehat{p}_t^d .

$$\text{PREREPR}_t^\bullet(tag) = \widehat{p}_t^d(tag) \cap \bigcap_{h \in \text{SubTerms}(t, tag)} \widehat{p}_t^c(h)$$

This memory pattern intersection may be empty, in which case representation generation fails and the input specification is deemed invalid since it provided incompatible memory patterns. Composability constraints are defined by iterating over subterms. Crucially, this yields a natural base case for constructors without subterms. $\text{REPRTY}_t^\bullet(tag)$ is a restriction of $\text{PREREPR}_t^\bullet(tag)$, where formulas are erased, yielding a memory type. We will now define each of these constraints.

5.1.2 Distinguishability Constraints

Distinguishability constraints $\widehat{p}_t^d(tag)$ encode that, for any given sum type t , the head constructor of any memory value can be identified by inspecting some predefined memory locations. Concretely, we accumulate all constraints on a given constructor, and then merge all constructor-specific constraints.

Let t such that $\Gamma(t) = \sum_n K_i(\overline{t_{i,j}^{n_i}})$ and K_i one of its constructors. We collect constraints by iterating over each discriminant and its constructors, collecting coarse-grained contexts, word operations, widths and values. For unsized memory holes (i.e., $*$), we use the maximal value width.

$$\mathcal{C} =_{\text{def}} \left\{ \begin{array}{l} (\widehat{h}, \ell, \text{ops}, z) \text{ where} \\ (\widehat{h}', \text{split}) \in \text{DISCRS}_t^\bullet \\ (K_i \mapsto z) \in \text{split} \\ \widehat{h}[\text{ops}_{\ell'}] = \widehat{h}' \\ \ell = (\text{if } (\ell' \text{ is } *) \text{ then } \max_{(K \mapsto z) \in \text{split}} \text{width}(z) \text{ else } \ell') \end{array} \right\}$$

We now define the *specific* distinguishability pattern $\widehat{p}_t^d(K_i)$ that applies to a given constructor K_i by instantiating contexts, intersecting the resulting patterns and taking the conjunction of word formulas:

$$\widehat{p}_t^d(K_i) = \bigcap_{(\widehat{h}_i, \ell_i, \text{ops}_i, z_i) \in \mathcal{C}} \widehat{h}_i[\text{Word}_{\ell_i}(\text{ops}_i = z_i)]$$

The *generic* distinguishability pattern of t captures the structural characteristics common to all constructors:

$$\widehat{p}_t^d(\top) = \bigcup_{K \in \text{Constrs}(t)} \widehat{p}_t^d(K)$$

5.1.3 Composability Constraints

Composability constraints $\widehat{p}_t^c(h)$ encode that the memory representation of a value that fits h must contain the memory representation of its subvalue at $\text{POSITION}_t^\bullet(h)$. Our goal is to collect representation constraints for each subterm position h , then merge them all.

Let t' be the subtype of t at the hole in h . Inductively, we have defined the *generic representation type* of t' : $\widehat{\tau}' = \text{REPRTY}_{t'}^\bullet(\top)$. Let $\widehat{h}[\text{ops}_{\ell_0}] = \text{POSITION}_t^\bullet(h)$. For the representation to be consistent, the memory hole at \widehat{h} must have type $\widehat{\tau}'$. There are two cases:

- The memory value at position \widehat{h} is a proper ℓ -bit wide word. We thus have $\widehat{\tau}' = \text{Word}_\ell$ and $\ell_0 = \ell$ or $\ell_0 = *$. To link the various constraints, we introduce x_h , a variable that contains the integer contents at position h . We output the constraint $\widehat{p}_t^c(h) = \widehat{h}[\text{Word}_\ell(\text{ops} = x_h)]$.

- The memory value at position \widehat{h} is a not a word (for instance, it is a block). We necessarily have $\underline{\text{ops}}_{\ell_0} = \blacksquare_{\ast}$. Since the stored memory value is not a word, we simply recall the constraints pertaining to the position h with its type t' using an application pattern: $\widehat{p}_t^c(h) = \widehat{h}[t'@h]$.

Example 8 (Running example: constraints). Consider the type t_{ABC} from our running example and the input ingredients of Example 6:

- t_{ABC} is a sum type with three variants. For the sake of brevity, we only give $\widehat{p}_{t_{\text{ABC}}}^d(C)$, which is an intersection of two memory patterns since two discriminants involve C :

$$\widehat{p}_{t_{\text{ABC}}}^d(C) = \text{Ptr}_{64,3} \{ \{ \text{Word}_{64} (\blacksquare = 2), \dots \} \cap \text{Word}_1 (\blacksquare[0,0] = 0)$$

This pattern matches memory values that have characteristics of both a 64-bit-wide, 3-bit-aligned pointer to a block with at least one field, whose first field is a 64-bit-wide word whose contents \blacksquare satisfy $\blacksquare = 2$, and of a word or address of at least one bit whose lowest bit is 0.

- t_{ABC} has one subterm $C(\square)$. Its composable pattern is:

$$\begin{aligned} \widehat{p}_{t_{\text{ABC}}}^c(C(\square)) &= \text{Ptr}_{64,3} \{ \{ -, \blacksquare_{\ast}, \dots \} [\text{Word}_{64} (\blacksquare = x_{C(\square)})] \\ &= \text{Ptr}_{64,3} \{ \{ -, \text{Word}_{64} (\blacksquare = x_{C(\square)}) \} \} \end{aligned}$$

which expresses that the representation of $C(n)$ must be a pointer to a block with at least two fields, whose second field is a word containing an (unknown) integer value y such that $y = x_{C(\square)}$, where $x_{C(\square)}$ denotes the integer contents of $\text{REPR}_{t_{\text{int}}}^{\bullet}(n)$.

◆

5.2 Representation Completion

We have now collected all constraints on our representation into a memory pattern that contains formulas. The next step of representation synthesis consists in *solving* each word formula. More precisely, given $\text{PREREPR}_t^{\bullet}(K)$ for some constructor K of t , we consider each memory context \widehat{h} such that the memory pattern $\text{Word}_{\ell_{\widehat{h}}}(\phi_{\widehat{h}})$ is at position \widehat{h} in the memory pattern $\text{PREREPR}_t^{\bullet}(K)$.

Let x_{h_1}, \dots, x_{h_m} be all the subvalue representation variables that appear in $\phi_{\widehat{h}}$. Let ℓ_1, \dots, ℓ_m the widths of these subvalues' representations (i.e., $\text{REPR}_{t_k}^{\bullet} = \text{Word}_{\ell_k}$). Our goal is to synthesise a function that, given the subvalues' representations, builds the word of the parent value (which we denoted \blacksquare in our formula). Let $y_{\widehat{h}}$ this function from bitvectors of widths ℓ_1, \dots, ℓ_m to an $\ell_{\widehat{h}}$ -bit-wide bitvector. Let $\phi'_{\widehat{h}} = \phi_{\widehat{h}}[\blacksquare \rightarrow y_{\widehat{h}}(x_{h_1}, \dots, x_{h_m})]$ the formula with \blacksquare substituted with its (symbolic) value. Solving $\phi'_{\widehat{h}}$ then consists in finding a function interpretation for y that satisfies it, denoted $f_{\widehat{h}}$. If $\phi'_{\widehat{h}}$ is unsatisfiable, the input specification is deemed invalid.

We solve each such formula $\phi'_{\widehat{h}}$ in $\text{PREREPR}_t^{\bullet}(K)$ and obtain suitable $f_{\widehat{h}}$ interpretations. This gives us fully set values for each memory word, allowing us to define the representation function for source values of constructor K . In order to define the full representation function, we perform the same procedure for each constructor K_i of t . We finally set:

$$\begin{aligned} \text{REPR}_t^{\bullet}(K_i(\overline{v_j})) &= \text{PREREPR}_t^{\bullet}(K_i) [\text{Word}_{\ell_{\widehat{h}}}(\phi_{\widehat{h}}) \rightarrow \text{Word}_{\ell_{\widehat{h}}}(e_{\widehat{h}})] \\ \text{where } e_{\widehat{h}} &= f_{\widehat{h}}(\text{REPR}_{t_1}^{\bullet}(v_{h_1}), \dots, \text{REPR}_{t_m}^{\bullet}(v_{h_m})) \end{aligned}$$

Example 9 (Running example: completion). We admit the following representation function for the primitive type t_{int} : $\text{REPR}_{t_{\text{int}}}^{\bullet}(n) = \text{Word}_{64}(2 \times n + 1)$ (see Section 6). Let us now build $\text{REPR}_{t_{\text{ABC}}}^{\bullet}(K)$ for each constructor $K \in \{A, B, C\}$:

- $K = A$. Collecting constraints yields this unsolved representation:

$$\begin{aligned} \text{PREREPR}_{t_{\text{ABC}}}^{\bullet}(A) &= \text{Word}_{64}(\blacksquare[0,0] = 1) \cap \text{Word}_{64}(\blacksquare = 1) \\ &= \text{Word}_{64}(\blacksquare[0,0] = 1 \wedge \blacksquare = 1) \end{aligned}$$

from which we extract the formula $\phi = (y[0,0] = 1 \wedge y = 1)$ whose unknown is $y \in \{0,1\}^{64}$ (a bitvector of size 64). The model $y = \text{Word}_{64}(1)$ satisfies ϕ , hence $\text{REPR}_{t_{\text{ABC}}}^{\bullet}(A) = \text{Word}_{64}(1)$. Similarly, we obtain $\text{REPR}_{t_{\text{ABC}}}^{\bullet}(B) = \text{Word}_{64}(3)$.

- $K = C$. After simplification, the constraints for C (from Example 8) give the following unsolved representation, in which the 64-bit-wide bitvector SMT variable $x_{C(\square)}$ symbolises the representation of the subvalue “under C ”:

$$\text{PREREPR}_{t_{\text{ABC}}}^{\bullet}(C) = \text{Ptr}_{64,3} \left\{ \left\{ \begin{array}{l} \text{Word}_{64}(\blacksquare = 2), \\ \text{Word}_{64}(\blacksquare = x_{C(\square)}) \end{array} \right\} \right\}$$

The first word description leads to a constant value of 2, the second one is solved by $y = x_{C(\square)}$ and we finally get:

$$\text{REPR}_{t_{\text{ABC}}}^{\bullet}(C(n)) = \text{Ptr}_{64,3} \left\{ \text{Word}_{64}(2), \text{Word}_{64}(\text{REPR}_{t_{\text{int}}}^{\bullet}(n)) \right\}$$

◆

5.3 Synthesis of Pattern Matching Primitives

We have previously asserted that the specification yield a consistent memory layout and synthesised a representation function to build memory values? We finally synthesise the pattern matching primitives introduced in Section 3.2.

UNWRAP_t[•] : *SubTerms*(t) → *Exprs*

$\text{UNWRAP}_t^{\bullet}(h)$ builds an expression to access the representation of the subvalue at position h in a value of type t . The input ingredient $\text{POSITION}_t^{\bullet}(h)$ returns a memory context \hat{h} providing this information. To obtain $\text{UNWRAP}_t^{\bullet}(h)$, we translate \hat{h} to a *target expression*, as defined in Fig. 5. For instance, if $\text{POSITION}_t^{\bullet}(h) = \left\{ \left\{ -, \dots, \square_i, \dots, - \right\} \right\}$, then $\text{UNWRAP}_t^{\bullet}(h) = \Delta.i$.

INSPECTANDSPLIT_t[•] : (*Constrs*(t) → *Trees*) → *Trees*

$\text{INSPECTANDSPLIT}_t^{\bullet}(\text{branch})$ builds a decision tree given a function *branch* from t constructors to decision trees. This primitive is more complex to build, and requires careful exploration of t discriminants to emit the minimal amount of switch nodes. We proceed by iteratively generating intermediate decision trees for decreasing sets of possible constructors, starting from all constructors of t . We then glue these intermediate decision trees with switch nodes.

More precisely, we define the decision tree $\mathcal{T}_{\mathcal{K}}$ by induction on the current set of constructors \mathcal{K} :

- If $\mathcal{K} = \{K\}$ contains a single constructor, we return $\mathcal{T}_{\mathcal{K}} = \text{branch}(K)$.
- Otherwise, let us consider H the set of discriminants which effectively distinguish values in \mathcal{K} :

$$H = \left\{ \hat{h} \mid \begin{array}{l} (\hat{h}, \text{split}) \in \text{DISCRS}_t^{\bullet} \\ \exists K, K' \in \mathcal{K} \cap \text{split}, \text{split}(K) \neq \text{split}(K') \end{array} \right\}$$

H might contain several discriminants. We choose the “shallowest” one, i.e. that doesn’t look far into the representation, and name it \hat{h}_0 . Thanks to the precondition on $\text{DISCRS}_t^{\bullet}$, we show (in Appendix C) that this shallowest discriminant always exists. By our distinguishability constraints, \hat{h}_0 is a valid context for all values whose constructor is in \mathcal{K} .

We can now build our switch node. Let e the target expression returning the value located at \hat{h}_0 (built similarly to UNWRAP^{\bullet}). By definition, the possible values of e are $V = \{z \mid K \in \mathcal{K} \cap \text{split}, \text{split}(K) = z\}$. For each z_i in V , let $\mathcal{K}_{z_i} = \{K \in \mathcal{K} \mid \text{split}(K) = z_i \vee K \notin \text{split}\}$ the constructors compatible with this result. Let $\mathcal{K}_{\top} = \{K \in \mathcal{K} \mid K \notin \text{split}\}$ the remaining constructors. We finally obtain:

$$\mathcal{T}_{\mathcal{K}} = \text{switch}(e) \left\{ \begin{array}{l} z_0 \mapsto \mathcal{T}_{\mathcal{K}_{z_0}} \\ \vdots \\ z_{m-1} \mapsto \mathcal{T}_{\mathcal{K}_{z_{m-1}}} \\ \top \mapsto \mathcal{T}_{\mathcal{K}_{\top}} \end{array} \right. \quad \text{where } z_i \in V$$

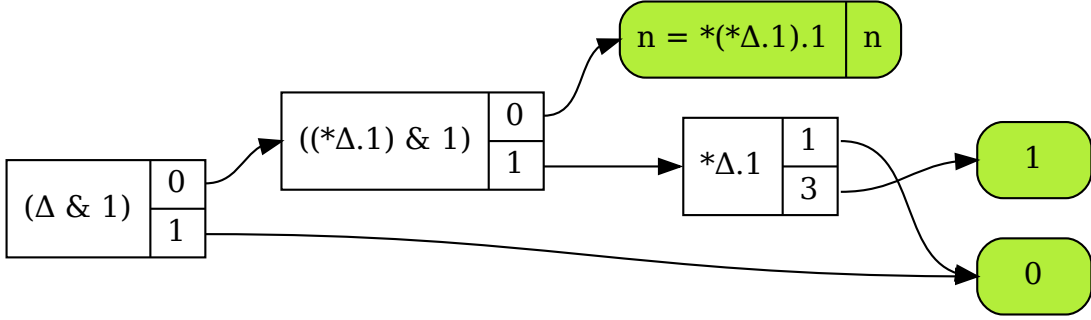


Figure 10: Output of `ribbit` for our running example $\text{COMPILE}_t(m_0)$ with the OCaml representation.

Appendix C shows the detailed algorithm and its proof.

Finally, our *representation-parametrised* COMPILE_t function presented in Section 3.2 is complete and can compile any pattern matching on t values to an equivalent decision tree. We state its correctness w.r.t. the memory layout:

Theorem 1. *Let Γ a type environment; $t \in \Gamma$, m a matching such that $\Gamma \vdash m : t \rightarrow \Gamma'$ and $\mathcal{T} = \text{COMPILE}_t(m)$. Then for any value v of type t such that $m \triangleright v \rightarrow i, \sigma$, the evaluation of \mathcal{T} on $\text{REPR}_t^\bullet(v)$ succeeds with $i, \{x \mapsto \text{REPR}_{\Gamma(x)}^\bullet \sigma(x)\}$.*

Example 10 (Running example: $\text{INSPECTANDSPLIT}^\bullet$ generation). Let us build $\mathcal{T}_{A,B,C}$. The shallowest discriminant (See Example 6) is $D_{A,B,C}$ that inspects the lowest bit of the representation, thus $e = \Delta \ \& \ 1$ that has values 0 (for C) or 1 (for A or B). The recursion for C gives $\text{branch}(C)$. The recursion using $D_{A,B}$ is similar. We finally obtain the helper function:

$$\text{INSPECTANDSPLIT}_{t_{ABC}}^\bullet(\text{branch}) =_{\text{def}} \left\{ \begin{array}{l} 0 \mapsto \text{branch}(C) \\ 1 \mapsto \text{switch}(\Delta) \left\{ \begin{array}{l} 1 \mapsto \text{branch}(A) \\ 3 \mapsto \text{branch}(B) \end{array} \right\} \end{array} \right\}$$

After generating all helper functions, a final call to $\text{COMPILE}_t(m_0)$ produces the final decision tree depicted in Fig. 10. Notice that switch nodes now use concrete expressions operating on memory values to extract subvalues' representations and distinguish between constructors.

For space reasons, we only gave a partial development of our running example in the OCaml representation. For the (far too extensive) details see Appendix D. \blacklozenge

6 Extensions

So far, we only showed the basic pattern constructs: sum types. Modern pattern languages have many additional constructs. We now sketch several extensions of our framework. These extensions are implemented in our prototype.

Reference and product types References and product are degenerate cases of sum types with only one case. Our framework trivially extends to these constructs, with the simplification that DISCRS^\bullet is always empty for these types (as there is nothing to be distinguished).

Primitive types Users may wish to pattern match on primitive types. At first glance, primitive types are similar to sum types where each constant is a unit constructor. Unfortunately, the number of distinct values (2^{64} for native integers) makes such definition of DISCRS^\bullet too large and the solving step intractable. Furthermore, primitive types are not always finite (big integers).

To alleviate this issue, we extend the notion of DISCRS^\bullet to work on *sets of values* which should be treated the same, as if we had several sum type constructors with a single field containing the individual value. Formally, for each T a primitive type, we require an additional ingredient VALSETS_T^\bullet which defines a partition of the values of T . We then require discriminants to work on these sets of values: $\text{DISCRS}_T^\bullet \subset \widehat{\text{Contexts}} \times (\text{VALSETS}_T^\bullet \rightarrow \mathbb{Z})$. The rest of our synthesis procedure readily extends to these additional specifications.

Bit-stealing A more interesting extension is an optimisation known as *bit-stealing* in which address alignment bits are used to store additional integer values. For instance, our RBT example can be represented even more compactly by using aligned pointers and storing the color in the lowest bit. Our framework already support leveraging pointer alignment: the OCaml representation exploits the fact that the lowest bit of pointers is always 0. To bit-steal, we extend pointers appropriately: $\text{Ptr}_{\ell,a}(z) \rightarrow \hat{v}$ is a ℓ -wide pointer to \hat{v} with z in the a lowest-bits. Pattern intersection and solving are easily adapted.

7 Related Works

7.1 Memory Representation and ADTs

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 optimised 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 behaviour and parallelism. Our work would augment these approaches with efficient decision tree generation.

7.2 Pattern Matching Compilation

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.

As for pattern-matching compilation, many works exists since the first introduction of pattern matrices in the context of the LML language Augustsson (1985). Fessant and Maranget (2001) first proposed optimisation 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. Most approaches to improve the current state-of-art algorithms rely on heuristics for the choice of column to split 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 has a much bigger impact on performance.

8 Conclusion

We have presented a specification and synthesis framework for *composable* memory representation, suitable for Algebraic Data Types implementation and compilation. Our framework not only synthesise optimising compilation procedure, but also checks the valid of the provided specification. As a case study, we used our method to recover a complete description for the OCaml representation and associated pattern-matching compilation. To our knowledge, this is the first generic specification of memory representation linked with Algebraic Data Types compilation purpose. 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 formalisation and description of new optimisation techniques for memory representation. In recent versions, Rust has been introducing more and more complex memory representation optimisations, which are so far unspecified and could benefit from our synthesis framework. Furthermore, we believe there is a trove of optimisations yet to be explored when it comes to the memory representation of Algebraic Data Types. Some promising leads are to apply super-optimisation to individual performance-sensitive data structures, or to allow programmers to specify whether types should be optimised for space, cache behaviour, or even sharing. We hope this work serves as a stepping stone for these further optimisations.

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
- 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>
- 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>
- Jacques Garrigue. 1998. Programming with polymorphic variants. In *ML Workshop*, Vol. 13. Baltimore.
- 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
- 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
- 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. 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>
- 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>

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

A Source patterns

In this appendix, we give an *exhaustive* presentation of our source pattern language, including its complete formal semantics. It includes all the extensions described in Section 6. This extends Section 3.1.

We first recall the complete grammar in Fig. 11.

Dynamic Semantics Fig. 12 defines the dynamic operational semantics of patterns. $p \triangleright v \rightarrow \sigma$ stands for “p matches value v and binds the variables in σ ”. Most of the rules are straightforward, with the following peculiarities:

- The bound value environment returned by the matching is populated by variables, through the VAR rule.
- We enforce that there is no shadowing: variables must be bound only once, as asserted by the side-condition in the 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 judgement in rule MATCHING: $\text{Match } \{p_1 \mid \dots \mid p_n\} \triangleright v \rightarrow i, \sigma$ which behaves as the normal judgement, 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.

Typing Fig. 13 defines the typing judgement $\Gamma \vdash p : t \rightarrow \Gamma'$: “pattern p has type t and binds variables whose types are defined in environment Γ' ”. Typing follows the semantics closely:

- As before, the bound typing environment Γ' 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 by the ALT rule.

B Memory patterns

Memory patterns are memory value skeletons in which possible values of address alignment bits and word contents are constrained by SMT formulae. They are notably used for structural interpolation, in Section 5.1. In this section, we flesh out their auxiliary operations, notably union, intersection and instantiation.

We first recall their grammar in Fig. 14. $t@h$ denotes the application of PREREPR_t^\bullet to the subvalue at position h . Note the extension of pointer pattern to accommodate for bit-stealing. Other cases closely follow the grammar of memory values: fixed-width words whose contents are specified by an SMT formula ϕ , fixed-width and fixed-alignment pointers (to another memory pattern) whose address alignment bits are specified by an SMT formula ϕ , blocks containing a finite number of fields (which are either memory patterns or the wildcard $_$) and disjoint unions of memory patterns. ϕ is either \top (always true) or a conjunction of clauses $\bigwedge_i (\text{ops}_i = x_i)$ where each x_i is either a concrete integer value or a subterm symbol and each ops_i is a sequence of operations applied to the word contents.

Union and intersection Memory patterns are equipped with intersection and union operations.

The intersection of two memory patterns, defined in Fig. 15, is empty if they are incompatible (that is, if they specify different memory structures) or another memory pattern that captures both patterns’ structural information and word contents otherwise. Intersection performs width extension: the intersection of two words is a word of the same width as the widest initial word. Also note that $(t@h)$ only intersects with memory patterns that do not constrain memory values more than $\text{PREREPR}_t^\bullet(\top)$ – indeed, $(t@h)$ denotes any memory value that matches $\text{PREREPR}_t^\bullet(\top)$.

The union of two memory patterns, defined in Fig. 16, merges their common memory structures and adds disjoint unions between incompatible structures and different SMT formulae. The resulting memory pattern captures all information that applies to any memory value that matches either of the initial patterns.

Instantiation *Memory context instantiation*, defined in Fig. 17, formalizes the substitutions in memory patterns. $\hat{h}[\hat{p}]$ is a memory pattern whose structure is that of the memory context \hat{h} , where the holes has been replaced by \hat{p} enhanced with the inner word operations.

C Decision tree generation algorithm for INSPECTANDSPLIT[•]

In Section 5.3, we give a definition of INSPECTANDSPLIT[•]_{*t*} using a recursive algorithm to build decision trees. This definition uses the concept of *shallowest* context, without defining it. In this appendix, we flesh out this definition and formalize and prove the algorithm to build decision trees.

C.1 Memory context relations

We first define a *partial prefix order* \prec on memory contexts in Fig. 18. Intuitively, if $\hat{h} \prec \hat{h}'$, then any memory value that is compatible with \hat{h}' is also compatible with \hat{h} . Accordingly, we assert that any transformation on words is a prefix of pointer dereferencing, as the former is compatible with any word or pointer of appropriate width whereas the latter only applies to pointers.

Two memory contexts may be “compatible” in that they both apply to at least one common memory value, without one being a prefix of the other. We formalise this relation in Fig. 19.

We also define a “semi-ordered compatibility” relation \bowtie from which we derive a notion of *minimum memory context*:

$$\hat{h} \bowtie \hat{h}' \iff \hat{h} = \hat{h}' \vee \hat{h} \prec \hat{h}' \vee \hat{h}' \prec \hat{h} \bowtie \hat{h}' \quad \min \hat{H} = \left\{ \hat{h} \in \hat{H} \mid \forall \hat{h}' \in \hat{H}, \hat{h} \bowtie \hat{h}' \right\}$$

Finally, the following results are needed to prove the existence of at least one minimal relevant discriminant:

Lemma 1. *Let \hat{H} a non-empty, finite set of memory contexts and \hat{P} a set of memory patterns instantiating all contexts of \hat{H} (i.e., $\forall \hat{h} \in \hat{H}, \exists \hat{p} \in \hat{P}, \exists \hat{p}', \hat{p} = \hat{h}[\hat{p}']$). Then we have*

$$\bigcap_{\hat{p} \in \hat{P}} \hat{p} \neq \emptyset \Rightarrow \min \hat{H} \neq \emptyset$$

Proof. We first show that for any memory contexts \hat{h} and \hat{h}' and any memory patterns \hat{p} and \hat{p}' such that $\hat{h}[\hat{p}] \cap \hat{h}'[\hat{p}'] \neq \emptyset$, we have $\hat{h} \bowtie \hat{h}'$ or $\hat{h}' \bowtie \hat{h}$. This result is immediate by induction on pattern intersection rules.

Therefore, for all $\hat{h}, \hat{h}' \in \hat{H}$, we have $\hat{h} \bowtie \hat{h}'$ or $\hat{h}' \bowtie \hat{h}$. Since \hat{H} is a finite ordered set, there exists at least one $\hat{h} \in \hat{H}$ such that $\forall \hat{h}' \in \hat{H}, \hat{h}' \not\prec \hat{h}$, that is $\forall \hat{h}' \in \hat{H}, \hat{h}' \bowtie \hat{h}$, which proves the existence of at least one minimal memory context of \hat{H} . \diamond

Lemma 2. *Let $\hat{h}_0, \hat{h}_1, \hat{h}_2$ and \hat{p}, \hat{p}' such that $\hat{h}_0 \bowtie \hat{h}_1$ and $\hat{h}_1[\hat{p}] \cap \hat{h}_2[\hat{p}'] \neq \emptyset$. Then we have $\hat{h}_0 \bowtie \hat{h}_2$.*

Proof. By induction on memory pattern intersection. \diamond

C.2 Decision tree generation algorithm

Recall the definition of INSPECTANDSPLIT[•]_{*t*}(*branch*) where *t* is a sum type and *branch* a function from *t* constructors to decision trees: INSPECTANDSPLIT[•]_{*t*}(*branch*) = $\mathcal{T}_{\text{Constrs}(t)}$, where $\mathcal{T}_{\mathcal{K}}$ is defined as follows for a set of constructors \mathcal{K} :

- if \mathcal{K} is empty, it is undefined and we omit any decision tree branch leading to \mathcal{T}_{\emptyset} ;
- if \mathcal{K} contains one constructor K , we branch to its associated continuation with $\mathcal{T}_{\{K\}} = \text{branch}(K)$;
- otherwise, we pick a *minimal relevant discriminant* $(\hat{h}, \text{split}) \in \text{DISCRS}_t^\bullet$ such that

$$\hat{h} \in \min \left\{ \hat{h} \mid \begin{array}{l} (\hat{h}, \text{split}) \in \text{DISCRS}_t^\bullet \\ \exists K, K' \in \mathcal{K} \cap \text{split}, \text{split}(K) \neq \text{split}(K') \end{array} \right\}$$

Let *e* the target expression returning the value located at \hat{h}_0 (built similarly to UNWRAP[•]). We build the following switch node:

$$\mathcal{T}_{\mathcal{K}} = \text{switch}(e) \left\{ c \mapsto \mathcal{T}_{\mathcal{K}(c)} \mid \begin{array}{l} c \in \mathbb{Z} \cup \{\top\} \\ \mathcal{K}_c \neq \emptyset \end{array} \right\}$$

with $\mathcal{K}_z = \{K \in \mathcal{K} \mid (K, z) \in \text{split} \vee K \notin \text{split}\}$ for $z \in \mathbb{Z}$ and $\mathcal{K}_{\top} = \{K \in \mathcal{K} \mid K \notin \text{split}\}$.

The correctness of this algorithm relies on the two following results.

Lemma 3 (Strictly decreasing constructor set). *Let $\mathcal{K} \subseteq \text{Constrs}(t)$ with $|\mathcal{K}| > 1$. For any source value v whose head constructor K is in \mathcal{K} , the execution of $\mathcal{T}_{\mathcal{K}}$ on $\text{REPR}_t^\bullet(v)$ branches to a recursively defined decision tree $\mathcal{T}_{\mathcal{K}'}$ with $K \in \mathcal{K}' \subsetneq \mathcal{K}$.*

Proof. Let $(\widehat{h}, \text{split})$ the minimal relevant discriminant picked for this step. Let v a source value whose head constructor K is in \mathcal{K} . Let us first remark that in any case, the execution of $\mathcal{T}_{\mathcal{K}}$ on $\text{REPR}_t^\bullet(v)$ branches to $\mathcal{T}_{\mathcal{K}(c)}$ for some $c \in \mathbb{Z} \cup \{\top\}$.

We first show that each such set $\mathcal{K}(c)$ is a strict subset of \mathcal{K} . By definition of $(\widehat{h}, \text{split})$, there are at least two distinct constructors K and K' in $\mathcal{K} \cap \text{split}$ such that $\text{split}(K) \neq \text{split}(K')$. It follows that at least one of these constructors is absent from each $\mathcal{K}(c)$ for $c \in \mathbb{Z} \cup \{\top\}$, hence $\mathcal{K}(c) \subsetneq \mathcal{K}$.

We now show that we always have $K \in \mathcal{K}(c)$. This is immediate if $K \notin \text{split}$. Otherwise, we have $K \in \text{split}$ and we show that the switch node branches to $\mathcal{T}_{\mathcal{K}(\text{split}(K))}$: upon evaluation on the input $\text{REPR}_t^\bullet(v)$, the switch expression e yields the memory value located at \widehat{h} . By construction, REPR_t^\bullet is such that this memory value is a word whose contents are $\text{split}(K)$. Indeed, for all $K \in \text{split}$, $\text{PREREPR}_t^\bullet(K)$ was built from an intersection of memory patterns including $\widehat{p} = \widehat{h}' [\text{Word}_\ell(\text{ops} = \text{split}(K))]$ where $\widehat{h}'[\text{ops}_{\ell'}] = \widehat{h}$ (ℓ and ℓ' are irrelevant here). This implies that the memory value located at \widehat{h}' in $\text{REPR}_t^\bullet(v)$ exists and is a word whose contents \blacksquare are such that $\text{ops} = \text{split}(K)$, that is, the memory value located at \widehat{h} in $\text{REPR}_t^\bullet(v)$ exists and is a word whose contents are equal to $\text{split}(K)$. \diamond

Lemma 4 (Existence of at least one minimal relevant discriminant). *Assume that the representation is valid, that is, $\text{REPR}_t^\bullet(K)$ exists for each constructor K of t . Then for all $\mathcal{K} \subseteq \text{Constrs}(t)$ such that $|\mathcal{K}| > 1$, we have*

$$\min \left\{ \widehat{h} \mid \begin{array}{l} (\widehat{h}, \text{split}) \in \text{DISCRS}_t^\bullet \\ \exists K, K' \in \mathcal{K} \cap \text{split}, \text{split}(K) \neq \text{split}(K') \end{array} \right\} \neq \emptyset$$

Proof. According to the distinguishability condition on DISCRS_t^\bullet , there exists a discriminant $(\widehat{h}, \text{split}) \in \text{DISCRS}_t^\bullet$ such that $K, K' \in \text{split}$ and $\text{split}(K) \neq \text{split}(K')$ for any constructor $K' \neq K$, therefore \widehat{H} is non-empty.

For any constructor K of type t , let

$$\widehat{H}_K = \widehat{H} \cap \left\{ \widehat{h} \mid \begin{array}{l} (\widehat{h}, \text{split}) \in \text{DISCRS}_t^\bullet \\ K \in \text{split} \end{array} \right\}$$

Recall the definition of the distinguishability pattern of any constructor K . Owing to representation validity, we have

$$\widehat{p}_t^d(K) = \bigcap_{(\widehat{h}_i, \ell_i, \text{ops}_i, z_i) \in \mathcal{C}} \widehat{h}_i [\text{Word}_{\ell_i}(\text{ops}_i = z_i)] \neq \emptyset$$

where

$$\mathcal{C} = \left\{ \begin{array}{l} (\widehat{h}, \ell, \text{ops}, z) \text{ where} \\ (\widehat{h}', \text{split}) \in \text{DISCRS}_t^\bullet \\ (K \mapsto z) \in \text{split} \\ \widehat{h}[\text{ops}_{\ell'}] = \widehat{h}' \\ \ell = (\text{if } (\ell' \text{ is } *) \text{ then } \max_{(K' \mapsto z') \in \text{split}} \text{width}(z) \text{ else } \ell') \end{array} \right\}$$

Notice that for any relevant discriminant $(\widehat{h}, \text{split}) \in \text{DISCRS}_t^\bullet$ such that $K \in \text{split}$, there exists $(\widehat{h}_i, \ell_i, \text{ops}_i, z_i) \in$

\mathcal{C} such that $\widehat{h} = \widehat{h}_i[\underline{\text{ops}}_i \ell'_i]$; we can write

$$\begin{aligned} & \widehat{h}_i[\text{Word}_{\ell_i}(\text{ops}_i = z_i)] \\ &= \widehat{h}_i[\underline{\text{ops}}_i \ell'_i[\text{Word}_{\ell_i}(\blacksquare = z_i)]] \\ &= \left(\widehat{h}_i[\underline{\text{ops}}_i \ell'_i]\right)[\text{Word}_{\ell_i}(\blacksquare = z_i)] \\ &= \widehat{h}[\text{Word}_{\ell_i}(\blacksquare = z_i)] \end{aligned}$$

That is, each considered memory context is instantiated once in this non-empty pattern intersection. We can thus apply Lemma 1 to show that for all $K \in \mathcal{K}$, $\min \widehat{H}_K \neq \emptyset$.

We must now prove that $\min \widehat{H} = \min \left\{ \min \widehat{H}_K \right\}_{K \in \mathcal{K}}$ is not empty. Let $K, K' \in \mathcal{K}$ such that $K \neq K'$. Using the same distinguishability condition argument as before, we have $\widehat{H}_K \cap \widehat{H}_{K'} \neq \emptyset$ and apply Lemma 2 on both constructors' distinguishability patterns to show that any minimal context for either constructor is compatible with all minimal contexts for the other constructor, and thus $\min \left(\min \widehat{H}_K \cup \min \widehat{H}_{K'} \right) \neq \emptyset$. \diamond

Using these two results, we show that $\text{INSPECTANDSPLIT}_t^\bullet$ descends through switch nodes that successively restrict the set of potential constructors, ending with the decision tree associated with the specific identified constructor.

We can thus assert the correctness of our pattern matching compilation scheme:

Theorem 2. *Let Γ a type environment; let $t \in \Gamma$ a type. Let m a matching such that $\Gamma \vdash m : t \rightarrow \Gamma'$ and $\mathcal{T} = \text{COMPILE}_t(m)$. Then for any value v of type t such that $m \triangleright v \rightarrow i, \sigma$, the evaluation of \mathcal{T} on $\text{REPR}_t^\bullet(v)$ succeeds with $i, \{x \mapsto \text{REPR}_{\Gamma(x)}^\bullet \sigma(x)\}$.*

D Full example: OCaml representation

Recall the type environment from our running example:

$$\begin{aligned} \Gamma &= \{ t \mapsto \text{None} + \text{Some}(t_{\text{ABC}}); \\ & \quad t_{\text{ABC}} \mapsto A + B + C(t_{\text{int}}); \\ & \quad t_{\text{int}} \mapsto u32 \} \end{aligned}$$

In this section, we detail the complete representation synthesis pipeline and pattern matching compilation for the OCaml representation, *with all extensions enabled*. (Note that the bit-stealing extension is not actually used in this representation, which is why all addresses' alignment bits are set to 0.)

D.1 Input specification

D.2 Primitive type t_{int}

In OCaml, all integer values follow the same memory layout. Let $S = [0; 2^{32} - 1]$ the set of $u32$ values. We encode any integer value $n \in S$ on 64 bits, then transform it through a one-bit left shift and an increment to set its pointer flag (lowest bit) to 1. As S is the only defined value set and no other representation constraints apply, no discriminant is needed for t_{int} . We thus define the following ingredients:

$$\text{VALSETS}_{t_{\text{int}}}^\bullet = \{S\} \quad \text{DISCRS}_{t_{\text{int}}}^\bullet = \emptyset \quad \text{POSITION}_{t_{\text{int}}}^\bullet(S(\square)) = \lfloor (\blacksquare - 1) / 2 \rfloor_{[0, 31]}_{64}$$

D.3 Sum types t_{ABC} and t

Sum types ingredients are similar across individual types. The main distinction done is between unit and non-unit constructors, which are represented by 64-bit-wide words and 64-bit-wide pointers to blocks, respectively.

We define the following ingredients for t_{ABC} :

$$\begin{aligned} \text{POSITION}_{t_{ABC}}^{\bullet}(C(\square)) &= \text{Ptr}_{64,3} \{ \{-, \blacksquare_{\ast}, \dots\} \} \\ \text{DISCRS}_{t_{ABC}}^{\bullet} &= \left\{ \begin{array}{l} D_{AB,C} = \left(\llbracket 0[\blacksquare, 0] \rrbracket_{\ast}, \left\{ \begin{array}{l} A \mapsto 1 \quad B \mapsto 1 \\ C \mapsto 0 \end{array} \right\} \right), \\ D_{A,B} = (\llbracket \blacksquare_{64} \rrbracket, \{A \mapsto 1 \quad B \mapsto 3\}), \\ D_C = (\text{Ptr}_{64,3} \{ \llbracket \blacksquare_{64}, \dots \rrbracket \}, \{C \mapsto 2\}) \end{array} \right\} \end{aligned}$$

We define the following ingredients for t :

$$\begin{aligned} \text{POSITION}_t^{\bullet}(\text{Some}(\square)) &= \text{Ptr}_{64,3} \{ \{-, \blacksquare_{\ast}, \dots\} \} \\ \text{DISCRS}_t^{\bullet} &= \left\{ \begin{array}{l} D_{None,Some} = \left(\llbracket 0[\blacksquare, 0] \rrbracket_{\ast}, \left\{ \begin{array}{l} None \mapsto 1 \\ Some \mapsto 0 \end{array} \right\} \right), \\ D_{None} = (\llbracket \blacksquare_{64} \rrbracket, \{None \mapsto 1\}), \\ D_{Some} = (\text{Ptr}_{64,3} \{ \llbracket \blacksquare_{64}, \dots \rrbracket \}, \{Some \mapsto 1\}) \end{array} \right\} \end{aligned}$$

D.4 Structural interpolation

Example 11 (Building the unsolved representation of a primitive type: t_{int}). t_{int} only has one value set S . No discriminant applies: by convention, we define $\widehat{p}_{t_{\text{int}}}^d(S) = _$. We derive its only composability pattern from $\text{POSITION}_{t_{\text{int}}}^{\bullet}(S(\square))$ to build its unsolved representation:

$$\begin{aligned} \text{PREREPR}_{t_{\text{int}}}^{\bullet}(\top) &= \text{PREREPR}_{t_{\text{int}}}^{\bullet}(S) \\ &= \widehat{p}_{t_{\text{int}}}^c(S) = \text{Word}_{64}(((\blacksquare - 1)/2)[0, 31] = x_{\square}) \\ \text{REPRTY}_{t_{\text{int}}}^{\bullet}(\top) &= \text{REPRTY}_{t_{\text{int}}}^{\bullet}(S) = \text{Word}_{64} \end{aligned}$$

◆

Example 12 (Composing word representations: t_{ABC}). We first compute the unsolved representations of both unit constructors A and B , which solely consists of their distinguishability pattern. Two discriminants apply: $D_{A,B}$ and $D_{AB,C}$. We decompose and recombine $D_{AB,C}$ into the memory pattern $\text{Word}_1(\blacksquare[0, 0] = 1)$ and get

$$\begin{aligned} \text{PREREPR}_{t_{ABC}}^{\bullet}(A) &= \widehat{p}_{t_{ABC}}^d(A) \\ &= \text{Word}_{64}(\blacksquare = 1) \cap \text{Word}_1(\blacksquare[0, 0] = 1) \\ &= \text{Word}_{64}(\blacksquare = 1 \wedge \blacksquare[0, 0] = 1) \\ \text{PREREPR}_{t_{ABC}}^{\bullet}(B) &= \widehat{p}_{t_{ABC}}^d(B) \\ &= \text{Word}_{64}(\blacksquare = 3) \cap \text{Word}_1(\blacksquare[0, 0] = 1) \\ &= \text{Word}_{64}(\blacksquare = 3 \wedge \blacksquare[0, 0] = 1) \end{aligned}$$

We therefore have $\text{REPRTY}_{t_{ABC}}^{\bullet}(A) = \text{REPRTY}_{t_{ABC}}^{\bullet}(B) = \text{Word}_{64}$.

For the non-unit constructor C , we need both a composability and a distinguishability memory pattern. Two discriminants apply: D_C and $D_{AB,C}$, hence

$$\begin{aligned} \widehat{p}_{t_{ABC}}^d(C) &= \text{Ptr}_{64,3}(\top) \rightarrow \{ \text{Word}_{64}(\blacksquare = 2), \dots \} \\ &\quad \cap \text{Word}_1(\blacksquare[0, 0] = 0) \\ &= \text{Ptr}_{64,3}(\top \wedge \blacksquare[0, 0] = 0) \rightarrow \{ \text{Word}_{64}(\blacksquare = 2), \dots \} \end{aligned}$$

We derive the composability pattern from $\text{POSITION}_{t_{ABC}}^{\bullet}(C(\square))$ and $\text{REPRTY}_{t_{\text{int}}}^{\bullet}(\top) = \text{REPRTY}_{t_{\text{int}}}^{\bullet}(S) = \text{Word}_{64}$:

$$\widehat{p}_{t_{ABC}}^c(C(\square)) = \text{Ptr}_{64,3}(\top) \rightarrow \{ \{-, \text{Word}_{64}(\blacksquare = x_{C(\square)}), \dots \} \}$$

Note that the representation of the subvalue “under C ” has been integrated as an SMT variable $x_{C(\square)}$. Hence

$$\begin{aligned} \text{PREREPR}_{t_{ABC}}^{\bullet}(C) &= \widehat{p}_{t_{ABC}}^d(C) \cap \widehat{p}_{t_{ABC}}^c(C(\square)) \\ &= \text{Ptr}_{64,3}(\top \wedge \blacksquare[0,0] = 0 \wedge \top) \\ &\quad \rightarrow \{\{\text{Word}_{64}(\blacksquare = 2), \text{Word}_{64}(\blacksquare = x_{C(\square)})\}\} \end{aligned}$$

and $\text{REPRTY}_{t_{ABC}}^{\bullet}(C) = \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{Word}_{64}\}\}$.

We finally compute the union of these partial representations to get the generic unsolved representation and its memory type:

$$\begin{aligned} &\text{PREREPR}_{t_{ABC}}^{\bullet}(\top) \\ &= \text{PREREPR}_{t_{ABC}}^{\bullet}(A) \cup \text{PREREPR}_{t_{ABC}}^{\bullet}(B) \cup \text{PREREPR}_{t_{ABC}}^{\bullet}(C) \\ &= \text{PREREPR}_{t_{ABC}}^{\bullet}(A) \sqcup \text{PREREPR}_{t_{ABC}}^{\bullet}(B) \sqcup \text{PREREPR}_{t_{ABC}}^{\bullet}(C) \\ &\quad \text{REPRTY}_{t_{ABC}}^{\bullet}(\top) \\ &= \text{REPRTY}_{t_{ABC}}^{\bullet}(A) \cup \text{REPRTY}_{t_{ABC}}^{\bullet}(B) \cup \text{REPRTY}_{t_{ABC}}^{\bullet}(C) \\ &= \text{Word}_{64} \cup \text{Word}_{64} \cup \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{Word}_{64}\}\} \\ &= \text{Word}_{64} \sqcup \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{Word}_{64}\}\} \end{aligned}$$

◆

Example 13 (Composing non-word representations: t). The unsolved representation of the unit constructor *None* is very similar to that of A and B from the previous example: two discriminants D_{None} and $D_{None,Some}$ apply and we get

$$\begin{aligned} \text{PREREPR}_t^{\bullet}(None) &= \widehat{p}_t^d(None) \\ &= \text{Word}_{64}(\blacksquare = 1) \cap \text{Word}_1(\blacksquare[0,0] = 1) \\ &= \text{Word}_{64}(\blacksquare = 1 \wedge \blacksquare[0,0] = 1) \\ \text{REPRTY}_t^{\bullet}(None) &= \text{Word}_{64} \end{aligned}$$

The unsolved representation of the non-unit constructor *Some* is the intersection of its only composability pattern and of its distinguishability pattern, involving two discriminants D_{Some} and $D_{None,Some}$.

$$\begin{aligned} \widehat{p}_t^d(Some) &= \text{Ptr}_{64,3}(\top) \rightarrow \{\{\text{Word}_{64}(\blacksquare = 1), \dots\}\} \\ &\quad \cap \text{Word}_1(\blacksquare[0,0] = 0) \\ &= \text{Ptr}_{64,3}(\top \wedge \blacksquare[0,0] = 0) \rightarrow \{\{\text{Word}_{64}(\blacksquare = 1), \dots\}\} \end{aligned}$$

Since $\text{REPRTY}_{t_{ABC}}^{\bullet}(\top)$ is not a word type, we integrate the subvalue representation as an application pattern rather than an SMT variable:

$$\widehat{p}_t^c(Some(\square)) = \text{Ptr}_{64,3}(\top) \rightarrow \{\{-, t_{ABC}@Some(\square), \dots\}\}$$

and we get

$$\begin{aligned} \text{PREREPR}_t^{\bullet}(Some) &= \widehat{p}_t^d(Some) \cap \widehat{p}_t^c(Some(\square)) \\ &= \text{Ptr}_{64,3}(\top \wedge \blacksquare[0,0] = 0) \rightarrow \{\{\text{Word}_{64}(\blacksquare = 1), \dots\}\} \\ &\quad \cap \text{Ptr}_{64,3}(\top) \rightarrow \{\{-, t_{ABC}@Some(\square), \dots\}\} \\ &= \text{Ptr}_{64,3}(\top \wedge \blacksquare[0,0] = 0 \wedge \top) \\ &\quad \rightarrow \{\{\text{Word}_{64}(\blacksquare = 1), t_{ABC}@Some(\square)\}\} \\ \text{REPRTY}_t^{\bullet}(Some) &= \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{REPRTY}_{t_{ABC}}^{\bullet}(\top)\}\} \end{aligned}$$

And we can finally build the generic unsolved representation and memory type:

$$\begin{aligned}
\text{PREREPR}_t^\bullet(\top) &= \text{PREREPR}_t^\bullet(\text{None}) \cup \text{PREREPR}_t^\bullet(\text{Some}) \\
&= \text{Word}_{64} (\blacksquare = 1 \wedge \blacksquare[0, 0] = 1) \\
&\sqcup \text{Ptr}_{64,3} (\top \wedge \blacksquare[0, 0] = 0 \wedge \top) \\
&\rightarrow \{\{\text{Word}_{64} (\blacksquare = 1), t_{\text{ABC}}@\text{Some}(\square)\}\} \\
\text{REPRTY}_t^\bullet(\top) &= \text{REPRTY}_t^\bullet(\text{None}) \cup \text{REPRTY}_t^\bullet(\text{Some}) \\
&= \text{Word}_{64} \sqcup \text{Ptr}_{64,3} \{\{\text{Word}_{64}, \text{REPRTY}_{t_{\text{ABC}}}^\bullet(\top)\}\}
\end{aligned}$$

◆

D.5 Representation completion

Example 14 (Representation completion of a primitive type: t_{int}). We complete the unsolved representation of the only value set of t_{int} . $\text{PREREPR}_{t_{\text{int}}}^\bullet(S) = \text{Word}_{64}(\phi)$ contains one formula at the position $\blacksquare_{\text{J}_*}$ that we express as the following atomic SMT problem:

$$\begin{aligned}
(\phi) y_{\blacksquare_{\text{J}_*}} &: \{0, 1\}^{32} \rightarrow \{0, 1\}^{64} \\
&\text{such that } \forall x_{\square} : \{0, 1\}^{32}, (y_{\blacksquare_{\text{J}_*}}(x_{\square}) - 1)/2[0, 31] = x_{\square}
\end{aligned}$$

where the variable x_{\square} symbolises the source integer value as a bitvector of width 32. We immediately solve it with the following model: $y_{\blacksquare_{\text{J}_*}}(x_{\square}) = (2 \times x_{\square} + 1)_{64}$, hence $\text{REPR}_{t_{\text{int}}}^\bullet(n) = \text{Word}_{64}(2 \times n + 1)$. ◆

Example 15 (Representation completion of t_{ABC}). Recall the unsolved representation patterns from the previous section. $\text{PREREPR}_{t_{\text{ABC}}}^\bullet(A) = \text{Word}_{64}(\phi)$ contains one formula at the position $\blacksquare_{\text{J}_*}$ that we express as the following atomic SMT problem:

$$(\phi) y_{\blacksquare_{\text{J}_*}} : \{0, 1\}^{64} \text{ such that } y_{\blacksquare_{\text{J}_*}} = 1 \wedge y_{\blacksquare_{\text{J}_*}}[0, 0] = 1$$

which we solve with the model $y_{\blacksquare_{\text{J}_*}} = (1)_{64}$, which satisfies ϕ , hence $\text{REPR}_{t_{\text{ABC}}}^\bullet(A) = \text{Word}_{64}(1)$.

We carry out the same steps to solve $\text{PREREPR}_{t_{\text{ABC}}}^\bullet(B)$: $y_{\blacksquare_{\text{J}_*}} = 3$ satisfies $y_{\blacksquare_{\text{J}_*}} = 3 \wedge y_{\blacksquare_{\text{J}_*}}[0, 0] = 1$ and thus $\text{REPR}_{t_{\text{ABC}}}^\bullet(B) = \text{Word}_{64}(3)$.

Solving $\text{PREREPR}_{t_{\text{ABC}}}^\bullet(C)$ is slightly more involved: it contains three formulas at positions $\widehat{h}_0 = \text{Ptr}_{64,3} \{\{\blacksquare_{\text{J}_*}, \dots\}\}$, $\widehat{h}_1 = \text{Ptr}_{64,3} \{\{-, \blacksquare_{\text{J}_*}, \dots\}\}$ and $\blacksquare_{\text{J}_*}$, which we express through the following atomic SMT problems:

$$\begin{aligned}
(\phi_0) y_{\widehat{h}_0} &\in \{0, 1\}^{64} \text{ such that } y_{\widehat{h}_0} = 2 \\
(\phi_1) y_{\widehat{h}_1} &\in \{0, 1\}^{64} \rightarrow \{0, 1\}^{64} \\
&\text{such that } \forall x_{C(\square)} : \{0, 1\}^{64}, y_{\widehat{h}_1}(x_{C(\square)}) = x_{C(\square)} \\
(\phi) y_{\blacksquare_{\text{J}_*}} &\in \{0, 1\}^3 \text{ such that } \top \wedge y_{\blacksquare_{\text{J}_*}}[0, 0] = 0 \wedge \top
\end{aligned}$$

We immediately solve these with the following models:

$$\begin{aligned}
y_{\widehat{h}_0} &= (2)_{64} \text{ satisfies } \phi_0 \\
y_{\widehat{h}_1}(x_{C(\square)}) &= x_{C(\square)} \text{ satisfies } \phi_1 \\
y_{\blacksquare_{\text{J}_*}} &= (0)_3 \text{ satisfies } \phi
\end{aligned}$$

and thus get, for any value v of type t_{int} :

$$\text{REPR}_{t_{\text{ABC}}}^\bullet(C(v)) = \text{Ptr}_{64,3}(0) \rightarrow \{\{\text{Word}_{64}(2), \text{Word}_{64}(\text{REPR}_{t_{\text{int}}}^\bullet(v))\}\}$$

◆

Example 16 (Representation completion of t). Recall the unsolved representation patterns from the previous section. $\text{PREREPR}_t^\bullet(\text{None}) = \text{Word}_{64}(\phi)$ contains one formula at the position \blacksquare_* that we express as the following atomic SMT problem:

$$(\phi) y_{\blacksquare_*} : \{0, 1\}^{64} \text{ such that } y_{\blacksquare_*} = 1 \wedge y_{\blacksquare_*}[0, 0] = 1$$

The model $y_{\blacksquare_*} = (1)_{64}$, satisfies ϕ , hence $\text{REPR}_t^\bullet(\text{None}) = \text{Word}_{64}(1)$.

$\text{PREREPR}_t^\bullet(\text{Some})$ contains three formulas at positions $\widehat{h}_0 = \text{Ptr}_{64,3} \{\{\blacksquare_*, \dots\}\}$, $\widehat{h}_1 = \text{Ptr}_{64,3} \{\{-, \blacksquare_*, \dots\}\}$ and \blacksquare_* that we express as the following atomic SMT problems:

$$\begin{aligned} (\phi_0) y_{\widehat{h}_0} &: \{0, 1\}^{64} \text{ such that } y_{\widehat{h}_0} = 1 \\ (\phi_1) y_{\widehat{h}_1} &: \{0, 1\}^{64} \rightarrow \{0, 1\}^{64} \text{ such that } \forall x : \{0, 1\}^{64}, y_{\widehat{h}_1}(x) = x \\ y_{\blacksquare_*} &: \{0, 1\}^3 \text{ such that } \top \end{aligned}$$

which we solve with the following models:

$$y_{\widehat{h}_0} = (1)_{64} \text{ satisfies } \phi_0 \quad y_{\widehat{h}_1}(x) = x \text{ satisfies } \phi_1 \quad y_{\blacksquare_*} = (0)_3 \text{ satisfies } \top$$

and we thus have, for any source value v of type t_{ABC} :

$$\begin{aligned} \text{REPR}_t^\bullet(\text{Some}(v)) &= \\ \text{Ptr}_{64,3}(0) &\rightarrow \{\{\text{Word}_{64}(1), \text{Word}_{64}(\text{REPR}_{t_{\text{ABC}}}^\bullet(v))\}\} \end{aligned}$$

Notice that $\text{REPR}_{t_{\text{ABC}}}^\bullet(v)$ corresponds to $x_{\text{Some}(\square)}$ in $\text{PREREPR}_t^\bullet(\text{Some})$. ◆

D.6 Pattern matching primitives

We finally synthesise pattern matching primitives UNWRAP^\bullet and $\text{INSPECTANDSPLIT}^\bullet$ for each sum type.

D.6.1 Composability primitives: UNWRAP^\bullet

Since the composability ingredients for the only subterm of both sum types t_{ABC} and t are identical, their composability primitives – which are generated by converting a memory context ingredient to a target expression – are also identical.

From the memory context

$$\text{POSITION}_{t_{\text{ABC}}}^\bullet(C(\square)) = \text{POSITION}_t^\bullet(\text{Some}(\square)) = \text{Ptr}_{64,3} \{\{-, \blacksquare_*, \dots\}\}$$

we obtain

$$\text{UNWRAP}_{t_{\text{ABC}}}^\bullet(C(\square)) = (*\Delta).1 \quad \text{UNWRAP}_t^\bullet(\text{Some}(\square)) = (*\Delta).1$$

With both types, the evaluation of this expression on the representation of a non-unit source value yields the representation of its subvalue under the constructor, as expected. Indeed, $*(\text{Ptr}_{64,3}(0) \rightarrow \{\{\widehat{v}_0, \widehat{v}_1\}\}).1$ evaluates to \widehat{v}_1 .

D.6.2 Distinguishability primitives: $\text{INSPECTANDSPLIT}^\bullet$

The different signatures of t_{ABC} and t result in distinct decision trees to distinguish between constructors. Indeed, $\text{INSPECTANDSPLIT}_{t_{\text{ABC}}}^\bullet$ features an extra switch node to distinguish between unit constructors A and B .

Example 17 (t_{ABC}). Let $\text{branch} : \{A, B, C\} \rightarrow \text{Trees}$. The generation of $\text{INSPECTANDSPLIT}_{t_{\text{ABC}}}^\bullet(\text{branch})$ begins with its toplevel decision tree $\mathcal{T}_{\{A, B, C\}}$.

The only minimal relevant discriminant for $\{A, B, C\}$ is $D_{A,B,C}$, which inspects the pointer flag to distinguish between unit and non-unit constructors. We emit the following switch node, in which recursively defined trees are present:

$$\mathcal{T}_{\{A,B,C\}} = \text{switch}(\Delta \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \mathcal{T}_{\{A,B\}} \\ 0 \mapsto \mathcal{T}_{\{C\}} \\ \top \mapsto \mathcal{T}_{\emptyset} \quad (\text{unreachable}) \end{array} \right\}$$

We now compute $\mathcal{T}_{\{A,B\}}$, whose minimal relevant discriminant is $D_{A,B}$, and $\mathcal{T}_{\{C\}}$:

$$\mathcal{T}_{\{A,B\}} = \text{switch}(\Delta) \left\{ \begin{array}{l} 1 \mapsto \mathcal{T}_{\{A\}} \\ 3 \mapsto \mathcal{T}_{\{B\}} \\ \top \mapsto F(\emptyset) \quad (\text{unreachable}) \end{array} \right\} \quad \mathcal{T}_{\{C\}} = \text{branch}(C)$$

Finally, we compute the last subtrees $\mathcal{T}_{\{A\}} = \text{branch}(A)$ and $\mathcal{T}_{\{B\}} = \text{branch}(B)$. Combining all previously defined subtrees yields the full decision tree: $\text{INSPECTANDSPLIT}_{t_{ABC}}^{\bullet}(\text{branch}) =$

$$\text{switch}(\Delta \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \text{switch}(\Delta) \left\{ \begin{array}{l} 1 \mapsto \text{branch}(A) \\ 3 \mapsto \text{branch}(B) \end{array} \right\} \\ 0 \mapsto \text{branch}(C) \end{array} \right\}$$

◆

Example 18 (*t*). Let $\text{branch} : \{\text{None}, \text{Some}\} \rightarrow \text{Trees}$. The generation of $\text{INSPECTANDSPLIT}_t^{\bullet}(\text{branch})$ begins with its toplevel decision tree $\mathcal{T}_{\{\text{None}, \text{Some}\}}$. Its only minimal relevant discriminant is $D_{\text{None}, \text{Some}}$, which inspects the pointer flag to distinguish between *None* and *Some*. This step is sufficient to fully distinguish all constructors: we have $\mathcal{T}_{\text{None}} = \text{branch}(\text{None})$, $\mathcal{T}_{\text{Some}} = \text{branch}(\text{Some})$ and obtain the following decision tree: $\text{INSPECTANDSPLIT}_t^{\bullet}(\text{branch}) =$

$$\text{switch}(\Delta \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \text{branch}(\text{None}) \\ 0 \mapsto \text{branch}(\text{Some}) \end{array} \right\}$$

◆

D.7 Pattern matching compilation

Recall the pattern matching problem from our running example:

$$m_0 = \text{Match} \left\{ \begin{array}{l} | (p_0) \ \text{None} \ | \ \text{Some}(A) \\ | (p_1) \ \text{Some}(B) \\ | (p_2) \ \text{Some}(C(n)) \end{array} \right\}$$

We can now apply our templated compilation procedure COMPILE_t to m_0 by instantiating $\text{INSPECTANDSPLIT}^{\bullet}$ and UNWRAP^{\bullet} calls.

Example 19 (OCaml pattern matching compilation of our running example). Our compilation procedure generates the following parametrised decision tree for the input m_0 :

$$\text{COMPILE}_t(m_0) = \text{INSPECTANDSPLIT}_t^{\bullet} \left(\begin{array}{l} \text{None} \mapsto \text{success}(0) \\ \text{Some} \mapsto \text{INSPECTANDSPLIT}_{t_{ABC}}^{\bullet} \left(\begin{array}{l} A \mapsto \text{success}(0) \\ B \mapsto \text{success}(1) \\ C \mapsto \text{success}(2, \{n \mapsto \text{UNWRAP}_{t_{ABC}}^{\bullet}(C(\square))\}) \end{array} \right) \\ [\Delta / \text{UNWRAP}_t^{\bullet}(\text{Some}(\square))] \end{array} \right)$$

We first determine whether we are dealing with *None* or a *Some* value; in the latter case, we determine whether the value “under *Some*” is a unit constructor or C , in which case we bind the representation of

the integer value “under C ” to the symbol n . We can now replace the symbolic calls to pattern matching primitives with their expressions to obtain the final decision tree:

$$\text{COMPILE}_t(m_0) = \text{switch}(\Delta \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \text{success}(0) \\ 0 \mapsto \text{switch}(((\ast\Delta).1) \ \& \ 1) \left\{ \begin{array}{l} 1 \mapsto \text{switch}((\ast\Delta).1) \left\{ \begin{array}{l} 1 \mapsto \text{success}(0) \\ 3 \mapsto \text{success}(1) \end{array} \right\} \\ 0 \mapsto \text{success}(2, \{n \mapsto ((\ast\Delta).1).1\}) \end{array} \right\} \end{array} \right\}$$

Notice how $\text{UNWRAP}_t^\bullet(\text{Some}(\square))$ calls are composed into the $\text{INSPECTANDSPLIT}_{t_{ABC}}^\bullet$ decision tree so as to inspect the representation of the *subvalue* “under *Some*”, which is $(\ast\Delta).1$. \blacklozenge

Patterns

$$\begin{array}{ll}
p ::= _ & \text{(Wildcard)} \\
| x & \text{(Variable)} \\
| (p \mid p') & \text{(Disjunction)} \\
| z \in \mathbb{Z} & \text{(Integer constant pattern)} \\
| \&p & \text{(Reference pattern)} \\
| \langle p_0, \dots, p_{n-1} \rangle & \text{(Product pattern)} \\
| K(p_0, \dots, p_{n-1}) & \text{(Constructor pattern)}
\end{array}$$
Matching

$$m ::= \text{Match } \{p_0 \mid \dots \mid p_{n-1}\} \quad \text{(Matching)}$$
Types

$$\begin{array}{ll}
t ::= T \text{ (e.g., i32, u16)} & \text{(Primitive type)} \\
| \&t & \text{(Reference type)} \\
| \prod_n t_i & \text{(Product type)} \\
| \sum_n K_i(t_{i,0}, \dots, t_{i,n_i-1}) & \text{(Sum type)} \\
\Gamma ::= \left\{ \begin{array}{l} x_0 : t_0; \dots; x_{n-1} : t_{n-1}; \\ t_0 \mapsto \tau_0; \dots; t_{n'-1} \mapsto \tau_{n'-1} \end{array} \right\} & \text{(Type environment)}
\end{array}$$
Values

$$\begin{array}{ll}
v ::= z \in \mathbb{Z} & \text{(Integer constant)} \\
| \&v & \text{(Reference)} \\
| \langle v_0, \dots, v_{n-1} \rangle & \text{(Product)} \\
| K(v_0, \dots, v_{n-1}) & \text{(Constructor)} \\
\sigma ::= \{x_0 \mapsto v_0; \dots; x_{n-1} \mapsto v_{n-1}\} & \text{(Value env.)}
\end{array}$$

Figure 11: PAT, our simplified language of patterns and types

$$\begin{array}{c}
\text{ANY} \\
\frac{}{_ \triangleright v \rightarrow \emptyset} \\
\\
\text{VAR} \\
\frac{}{x \triangleright v \rightarrow \{x \mapsto v\}} \\
\\
\text{CONSTANT} \\
\frac{}{z \triangleright z \rightarrow \emptyset} \\
\\
\text{REFERENCE} \\
\frac{p \triangleright v \rightarrow \sigma}{\&p \triangleright \&v \rightarrow \sigma} \\
\\
\text{TUPLE} \\
\frac{\forall i, p_i \triangleright v_i \rightarrow \sigma_i \quad \forall j \neq i, \sigma_i \cap \sigma_j = \emptyset}{\langle \bar{p}_i \rangle \triangleright \langle \bar{v}_i \rangle \rightarrow \bigcup \sigma_i} \\
\\
\text{CONSTRUCTOR} \\
\frac{\forall i, p_i \triangleright v_i, \sigma_i \quad \forall j \neq i, \sigma_i \cap \sigma_j = \emptyset}{K(\bar{p}_i) \triangleright K(\bar{v}_i) \rightarrow \bigcup \sigma_i} \\
\\
\text{ALTL} \\
\frac{p_1 \triangleright v \rightarrow \sigma}{(p_1 \mid p_2) \triangleright v \rightarrow \sigma} \\
\\
\text{ALTR} \\
\frac{p_1 \not\triangleright v \quad p_2 \triangleright v \rightarrow \sigma}{(p_1 \mid p_2) \triangleright v \rightarrow \sigma} \\
\\
\text{MATCHING} \\
\frac{p_i \triangleright v \rightarrow \sigma \quad \forall j < i, p_j \not\triangleright v}{\text{Match } \{p_1 \mid \dots \mid p_n\} \triangleright v \rightarrow i, \sigma}
\end{array}$$
Figure 12: Semantics of patterns: $p \triangleright v \rightarrow \sigma$

$$\begin{array}{c}
\text{ANY} \\
\Gamma \vdash _ : t \rightarrow \emptyset \\
\\
\text{VAR} \\
\Gamma \vdash x : t \rightarrow \{x : t\} \\
\\
\text{REFERENCE} \\
\frac{\Gamma \vdash p : t \rightarrow \Gamma'}{\Gamma \vdash \&p : \&t \rightarrow \Gamma'} \\
\\
\text{CONSTANT} \\
\frac{z \in T}{\Gamma \vdash z : T \rightarrow \emptyset} \\
\\
\text{TUPLE} \\
\frac{\forall i; \Gamma \vdash p_i : t_i \rightarrow \Gamma_i \quad \forall i, j; \Gamma_i \cap \Gamma_j = \emptyset}{\Gamma \vdash \langle \overline{p_i} \rangle : \prod t_i \rightarrow \bigcup \Gamma_i} \\
\\
\text{CONSTRUCTOR} \\
\frac{\exists i_0, K = K_{i_0} \quad \forall j; \Gamma \vdash p_j : t_{i,j} \rightarrow \Gamma_j \quad \forall i, j; \Gamma_i \cap \Gamma_j = \emptyset}{\Gamma \vdash K(\overline{p_j}) : \sum K_i(t_{i,j}) \rightarrow \bigcup \Gamma_j} \\
\\
\text{ALT} \\
\frac{\forall i; \Gamma \vdash p_i : t \rightarrow \Gamma'}{\Gamma \vdash p_1 \mid p_2 : t \rightarrow \Gamma'} \\
\\
\text{MATCHING} \\
\frac{\forall i; \Gamma \vdash p_i : t \rightarrow \Gamma'}{\Gamma \vdash \text{Match} \{p_1 \mid \dots \mid p_n\} : t \rightarrow \Gamma'}
\end{array}$$

Figure 13: Typing of patterns: $\vdash p : t \rightarrow \Gamma$

$$\begin{aligned}
\phi &::= \top \mid (\text{ops} = x_h) \mid (\text{ops} = z) \mid \phi \wedge \phi \\
\hat{p} &::= t@h \mid \text{Word}_\ell(\phi) \mid \text{Ptr}_{\ell,a}(\phi) \rightarrow \hat{p} \mid \{\{\overline{\hat{p}^?}\}\} \mid \bigsqcup \hat{p}
\end{aligned}$$

Figure 14: Memory patterns values and types

$$\begin{array}{c}
(t@h) \cap (t@h) = t@h \quad \frac{\text{PREREPR}_t^\bullet(\top) \cap \hat{p} = \text{PREREPR}_t^\bullet(\top)}{(t@h) \cap \hat{p} = t@h} \quad \frac{\ell \leq \ell'}{\text{Word}_\ell(\phi) \cap \text{Word}_{\ell'}(\phi') = \text{Word}_{\ell'}(\phi \wedge \phi')} \\
\\
\frac{\ell' \leq a}{\text{Ptr}_{\ell,a}(\phi) \rightarrow \hat{p} \cap \text{Word}_{\ell'}(\phi') = \text{Ptr}_{\ell,a}(\phi \wedge \phi') \rightarrow \hat{p}} \quad \frac{\hat{p} \cap \hat{p}' = \hat{p}''}{\text{Ptr}_{\ell,a}(\phi) \rightarrow \hat{p} \cap \text{Ptr}_{\ell,a}(\phi') \rightarrow \hat{p}'' = \text{Ptr}_{\ell,a}(\phi \wedge \phi') \rightarrow \hat{p}''} \\
\\
\frac{n \leq n' \quad \forall i < n, \hat{p}_i \cap \hat{p}'_i = \hat{p}''_i}{\{\{\overline{\hat{p}_i^n}, \dots\}\} \cap \{\{\overline{\hat{p}'_i^{n'}}, \dots\}\} = \{\{\overline{\hat{p}''_i}, \overline{\hat{p}''_{n+i}{}^{n'-n}}, \dots\}\}} \quad \frac{\forall i, \hat{p} \cap \hat{p}_i = \hat{p}'_i}{\hat{p} \cap \bigsqcup_n \hat{p}_i = \bigsqcup_n \hat{p}'_i}
\end{array}$$

Figure 15: Intersection

$$\begin{array}{c}
(t@h) \cup (t@h) = t@h \quad \frac{\text{PREREPR}_t^\bullet(\top) \cup \hat{p} = \text{PREREPR}_t^\bullet(\top)}{(t@h) \cup \hat{p} = t@h} \quad \text{Word}_\ell(\phi) \cup \text{Word}_\ell(\phi) = \text{Word}_\ell(\phi) \\
\\
\frac{\hat{p} \cup \hat{p}' = \hat{p}''}{\text{Ptr}_{\ell,a}\hat{p}(\phi) \cup \text{Ptr}_{\ell,a}\hat{p}'(\phi) = \text{Ptr}_{\ell,a}\hat{p}''(\phi)} \\
\\
\frac{\hat{p}_{i_0} \cup \hat{p}'_{i_0} = \hat{p}''_{i_0}}{\{\{\overline{\hat{p}_i^n}, \dots\}\} \cup \{\{\overline{\hat{p}_i^{i_0}}, \hat{p}'_{i_0}, \overline{\hat{p}_i^{n-i_0-1}}, \dots\}\} = \{\{\overline{\hat{p}_i^{i_0}}, \hat{p}'_{i_0}, \overline{\hat{p}_i^{n-i_0-1}}, \dots\}\}} \quad \frac{\forall i, \hat{p} \cup \hat{p}_i = \hat{p}'_i}{\hat{p} \cup \bigsqcup_n \hat{p}_i = \bigsqcup_n \hat{p}'_i} \\
\\
\frac{\text{no other rule applies}}{\hat{p} \cup \hat{p}' = \hat{p} \sqcup \hat{p}'}
\end{array}$$

Figure 16: Union

$$\begin{array}{c}
\blacksquare_*[\widehat{p}] = \widehat{p} \\
\frac{\text{ops}[\phi] = \phi'}{\underline{\text{ops}}_*[\text{Word}_\ell(\phi)] = \text{Word}_\ell(\phi') \quad \underline{\text{ops}}_\ell[\text{Word}_\ell(\phi)] = \text{Word}_\ell(\phi')} \\
\frac{\text{ops}[\phi] = \phi'}{\underline{\text{ops}}_*[\text{Ptr}_{\ell,a}(\phi) \rightarrow \widehat{p}] = \text{Ptr}_{\ell,a}(\phi') \rightarrow \widehat{p} \quad \underline{\text{ops}}_{a\ell}[\text{Ptr}_{\ell,a}(\phi) \rightarrow \widehat{p}] = \text{Ptr}_{\ell,a}(\phi') \rightarrow \widehat{p}} \\
\frac{\widehat{h}[\widehat{p}] = \widehat{p}'}{(\text{Ptr}_{\ell,a}\widehat{h})[\widehat{p}] = \text{Ptr}_{\ell,a}(\top) \rightarrow \widehat{p}'} \quad \frac{\widehat{h}[\widehat{p}] = \widehat{p}'}{\left\{ \{-, \dots, \widehat{h}_i, \dots, -\} \right\}[\widehat{p}] = \left\{ \{-, \dots, \widehat{p}'_i, \dots, -\} \right\}}
\end{array}$$

Figure 17: Memory context instantiation – $\text{ops}[\phi]$ substitutes each clause $\text{ops}' = x$ in ϕ with $\text{ops} \circ \text{ops}' = x$

$$\begin{array}{c}
\blacksquare_* \prec \underline{\text{op}}(\text{ops})_* \quad \blacksquare_\ell \prec \underline{\text{op}}(\text{ops})_\ell \quad \frac{\underline{\text{ops}}_* \prec \underline{\text{ops}}'_*}{\underline{\text{op}}(\text{ops})_* \prec \underline{\text{op}}(\text{ops}')_*} \quad \frac{\underline{\text{ops}}_\ell \prec \underline{\text{ops}}'_\ell}{\underline{\text{op}}(\text{ops})_\ell \prec \underline{\text{op}}(\text{ops}')_\ell} \\
\underline{\text{ops}}_* \prec \text{Ptr}_{\ell,a}\widehat{h} \quad \blacksquare_\ell \prec \text{Ptr}_{\ell,a}\widehat{h} \quad \underline{\text{ops}}_{a\ell} \prec \text{Ptr}_{\ell,a}\widehat{h} \quad \blacksquare_* \prec \left\{ \{-, \dots, \widehat{h}_i, \dots, -\} \right\} \\
\frac{\widehat{h} \prec \widehat{h}'}{\text{Ptr}_{\ell,a}\widehat{h} \prec \text{Ptr}_{\ell,a}\widehat{h}' \quad \left\{ \{-, \dots, \widehat{h}_i, \dots, -\} \right\} \prec \left\{ \{-, \dots, \widehat{h}'_i, \dots, -\} \right\}}
\end{array}$$

Figure 18: Partial prefix order

$$\begin{array}{c}
\frac{\text{op} \neq \text{op}'}{\underline{\text{op}}(\text{ops})_* \not\bowtie \underline{\text{op}}'(\text{ops}')_* \quad \underline{\text{op}}(\text{ops})_\ell \not\bowtie \underline{\text{op}}'(\text{ops}')_\ell} \quad \frac{\underline{\text{ops}}_* \not\bowtie \underline{\text{ops}}'_*}{\underline{\text{op}}(\text{ops})_* \not\bowtie \underline{\text{op}}(\text{ops}')_*} \\
\frac{\underline{\text{ops}}_\ell \not\bowtie \underline{\text{ops}}'_\ell}{\underline{\text{op}}(\text{ops})_\ell \not\bowtie \underline{\text{op}}(\text{ops}')_\ell} \quad \frac{\underline{\text{ops}}_* \not\prec \text{Ptr}_{\ell,a}\widehat{h}}{\underline{\text{ops}}_* \not\bowtie \text{Ptr}_{\ell,a}\widehat{h}} \quad \frac{\underline{\text{ops}}_{a\ell} \not\prec \text{Ptr}_{\ell,a}\widehat{h}}{\underline{\text{ops}}_{a\ell} \not\bowtie \text{Ptr}_{\ell,a}\widehat{h}} \\
\frac{i \neq i'}{\left\{ \{-, \dots, \widehat{h}_i, \dots, -\} \right\} \not\bowtie \left\{ \{-, \dots, \widehat{h}'_{i'}, \dots, -\} \right\}} \quad \frac{\widehat{h} \not\bowtie \widehat{h}'}{\text{Ptr}_{\ell,a}\widehat{h} \not\bowtie \text{Ptr}_{\ell,a}\widehat{h}' \quad \left\{ \{-, \dots, \widehat{h}_i, \dots, -\} \right\} \not\bowtie \left\{ \{-, \dots, \widehat{h}'_{i'}, \dots, -\} \right\}}
\end{array}$$

Figure 19: Non-comparable compatibility

The Inria logo is written in a red, cursive script font. It is centered within a white rounded rectangular box that has a subtle drop shadow.

**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

ISSN 0249-6399