

Elaborating Inductive Types to Custom Data Structures

AUTHOR(S), University of St Andrews, UK

ACM Reference Format:

Author(s). 2024. Elaborating Inductive Types to Custom Data Structures. 1, 1 (May 2024), 9 pages. <https://doi.org/10.1145/nnnnnnnn>. nnnnnnnn

1 INTRODUCTION

Functional languages offer a high degree of expressiveness using a very small set of primitives: most employ a variant of the typed lambda calculus, commonly with schemas for defining recursive or inductive types. Such constructs allow programmers to define data and logic in a natural, algebraic way that is amenable to abstraction. Data structures such as lists, trees, and even natural numbers are the archetypical examples of inductively defined data, which offer pattern-matching as a principled way to introduce ‘complete’ branching points in a program. Despite these advantages, functional languages generally do not provide programmers with many tools to specify how their programs should be represented on physical computers. For inductively-defined data, there is a fixed representation that the majority of functional languages utilise, which is to represent each constructor as a heap cell, and link chains of constructors together using pointer indirection. As a result, a list as an inductively defined data structure is stored as a linked list, and a natural number is stored as a unary number where each digit is an ‘empty’ heap cell. To get around this issue, most real-world implementations expose underlying machine primitive types such as contiguous arrays and bitvectors, and programmers are able to utilise these instead of the ‘algebraic’ inductively-defined types to make their functional programs more performant. In Haskell, the array package comes to mind, as well as the **Integer** and **Natural** primitives which utilise a GMP-style big-integer implementation internally.

1.1 The ‘Nat-hack’

In functional languages with dependent types, such as Idris or Lean, the issue of type representation is further complicated by the fact that inductive proofs come into the mix. Inductively defined types enjoy induction principles which can prove facts about them on a case-by-case basis. The standard example is the induction over the natural numbers **data Nat = Z | S Nat** given by

```
induction : (P : Nat -> Type) -> P Z -> ((m : Nat) -> P m -> P (S m)) -> (n : Nat) -> P n
```

This can be proven trivially by well-founded recursion and case analysis over **Z** and **S**. However, if natural numbers are not defined inductively, but rather opaquely as an intricate data structure like Haskell’s **Natural**, induction is no longer given ‘for free’, and must be manually proven internally in the language. On the other hand, if natural numbers are defined inductively, then many operations become inexcusably slow, such as multiplication

Author’s address: Author(s), University of St Andrews, St Andrews, UK, kt81@st-andrews.ac.uk.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

```

53      mul : Nat -> Nat -> Nat
54      mul Z b = Z
55      mul (S a) b = add (mul a b) b

```

taking $a \times b$ steps to compute for input numbers a and b . To solve this problem, Idris, Agda and Lean ‘short-circuit’ the default inductive type representation for natural numbers, to use arbitrarily-sized big integers for calculations whose digits are bitvectors, rather than unary numbers.

To describe how this trick works, assume we have access to a type **BigUInt** for arbitrarily-sized big integers, along with some primitive operations

```

63      bigZero, bigOne : BigUInt
64      bigAdd, bigMul, bigSub : BigUInt -> BigUInt -> BigUInt
65      bigIsZero : BigUInt -> Bool

```

In a well-typed input program, all occurrences of the zero constructor **Z** : **Nat** should be replaced with the constant **bigZero**, and all occurrences of the successor constructor **S** : **Nat** -> **Nat** should be replaced with the expression $\backslash x \rightarrow \text{bigAdd } x \text{ bigOne}$. For case analysis, each pattern matching expression **case** x **of** **Z** -> b | **S** n -> r n should be replaced with the conditional expression **if** **bigIsZero** x **then** b **else** r (**bigSub** x **bigOne**). Additionally, some basic functions on natural numbers should be replaced with more performant variants. The recursively-defined addition function **add** : **Nat** -> **Nat** -> **Nat** should be replaced with **bigAdd** and similarly **mul** should be replaced with **bigMul**. This way, the high-level program still appears to be using the inductive definition **Nat**, but upon compilation it uses **BigUInt** for efficient execution.

1.2 Beyond natural numbers

There are arguably many inductive types which could admit a more optimised representation than the default linked tree. The first which comes to mind is the type of lists **data List** $t = \text{Nil} \mid \text{Cons } t \text{ (List } t)$. There exist many representations of lists in memory, including flattened contiguous heap-backed arrays with dynamic resizing, singly-linked lists, doubly-linked lists, their circular variants, tree-based representations like binary search trees, their balanced variants, B and B+ trees, and segment trees. Each representation offers different performance characteristics for common operations such as appending, insertion, deletion, splicing, concatenating, and so on. Nevertheless, all of them are **List** in spirit; there is a ‘canonical’ bijection between a list and any of the aforementioned data structures. A functional program using the algebraic **List** type could potentially benefit from a different representation depending on the exact operations it performs and in what proportion. Not only lists, but structures such as trees, finite sets, as well as refinement predicates on types such as element proofs or parity proofs, could all be subject to more optimal representations. Since, at first glance, such an alteration could be done purely mechanically in a similar way to the ‘Nat-hack’, it is natural to wonder if this technique readily generalises to user-defined inductive types such that the transformation itself is specified in the same language.

1.3 Contributions

This paper develops an extension of dependently-typed lambda calculus with inductive constructions, in order to support the definition of custom representations of the inductive constructions, along with the specialisation of functions for fine tuning to the chosen representations. This system features correct-by-construction representations, awarding the programmer with a bijection proof between an inductive type and its representation, as well as an elaboration into a lower language with a guarantee that no inductive constructors or eliminators remain. The standard linked tree

representation of inductive types that is the hardcoded default in most implementations is incarnated as ‘just another representation’ in this system, special only because it can apply to *any* inductive type. The order of these developments in the paper are as follows:

- In section 2.1, the language λ_{PRIM} is introduced, which is a staged language with dependent types, whose object-level fragment contains some machine primitives that act as building blocks of data representations.
- In section 2.2, the language λ_{IND} is introduced as an extension of λ_{PRIM} , which allows the familiar definition of inductive types, living in a universe of ‘codes’ for object-level types.
- In section 2.3 The language λ_{REP} is introduced as the completion of λ_{IND} , containing a ‘representation’ construct that assigns concrete object-level codes to inductive types.
- In section 3, an elaboration procedure $\mathcal{R} : \lambda_{\text{REP}} \rightarrow \lambda_{\text{PRIM}}$ is formulated which eliminates inductive constructs through the defined representations, yielding the final program that can be staged and compiled.
- In section 4 The standard linked tree representation is recovered in λ_{REP} for any inductive type, and shown to be coherent.
- In section 5, the ‘Nat-hack’ is defined internally in λ_{REP} and shown to be coherent up to some notable assumptions. The system is further explored, showcasing representations for other inductive types.
- Finally, in section 6, some desirable extensions as part of future work are discussed.

2 A TYPE SYSTEM FOR DATA REPRESENTATIONS

In this section, we describe a type system for data representations in a staged language. We start by defining a core staged language λ_{PRIM} with Σ/Π types, identity types, and a universe of types \mathcal{U}_i for each stage i , as well as a set of object-level machine primitives. We then introduce inductive constructions in the meta-fragment of λ_{PRIM} to form λ_{IND} . Finally, we introduce data representations in λ_{IND} and extend the system with rules for data representations in λ_{REP} . Since staged dependent type systems and inductive constructions are well understood, we will focus on the novel aspects of data representations. The languages above form an inclusion hierarchy

$$\lambda_{\text{PRIM}} \subset \lambda_{\text{IND}} \subset \lambda_{\text{REP}}$$

and our goal is to describe a transformation from λ_{REP} to λ_{PRIM} , elaborating away inductive definitions with the help of data representations. Finally, the resulting program in λ_{PRIM} can be staged to the purely object-level language that represents the target architecture. When defining these languages, we will use a BNF-like syntax for terms and contexts, natural deduction-style typing rules, as well as CWF-style notation for contexts, types, and terms.

2.1 A core staged language with machine primitives, λ_{PRIM}

As a first step towards a type system for data representations, we informally describe a staged dependent type system λ_{PRIM} with Σ/Π types, identity types, and a universe of types \mathcal{U}_i for each stage $i \in \{0, 1\}$. This serves as a background system on which we introduce inductive constructions and data representations in. The raw syntax of λ_{PRIM} is given by the BNF grammar

$$t ::= \mathcal{U}_i \mid x \mapsto t \mid t \mid x \mid (x : t) \rightarrow t \mid (x : t) \times t \mid a =_A b \mid \pi_1 t \mid \pi_2 t \mid \mathbf{refl} \mid \uparrow t \mid \langle t \rangle \mid \sim t$$

This follows the standard typing rules of 2LTT with Σ , Π and identity types [3], without regard for the universe hierarchy, as this is orthogonal to the main focus of this work. Notably, $\mathcal{U}_0 : \mathcal{U}_0$ is the object-level universe, and $\mathcal{U}_1 : \mathcal{U}_1$ is the meta-level universe. We will use **let** _{t} notation for binding, though this is just syntactic sugar for redexes.

We need to add explicit substitutions here!

On top of this base syntax, we assume a certain set of primitives that exist in the object language. These will be used to form data representations. They are provided by the target architecture, which is represented by the object language:

- A type of booleans, $\text{Bool} : \mathcal{U}_0$, with constants $\text{true} : \text{Bool}$ and $\text{false} : \text{Bool}$, and operations and , or , and not . Elimination for booleans is also provided, in the form of

$$\text{ifthenelse} : (b : \text{Bool}) \rightarrow ((b =_{\text{Bool}} \text{true}) \rightarrow A) \rightarrow ((b =_{\text{Bool}} \text{false}) \rightarrow A) \rightarrow A.$$

- A type of machine words, $\text{Word} : \mathcal{U}_0$ with constants $0, 1, 2, \dots$ and standard binary numeric operations add , sub , mul , as well as Bool-valued comparison operations eq , lt , and gt . We will use the notation Word_n to denote the type $(w : \text{Word}) \times \text{lt } w \text{ } n =_{\text{Bool}} \text{true}$.
- A type of n -sized sequences of type A , $[A; n] : \mathcal{U}_0$, where $n : \text{Word}$ and $A : \mathcal{U}_0$. Sequences come with indexing operations $\text{get} : [A; n] \rightarrow \text{Word}_n \rightarrow A$ and $\text{set} : [A; n] \rightarrow \text{Word}_n \rightarrow A \rightarrow [A; n]$, as well as an initialisation operation $\text{repeat} : (a : A) \rightarrow (n : \text{Word}) \rightarrow [A; n]$. Sequences should be thought of as unboxed arrays, living on the stack.
- A boxing type constructor $\Box A : \mathcal{U}_0$ where $A : \mathcal{U}_0$, with boxing and unboxing operations $\text{box } a : \Box A$ and $\text{unbox } b : A$. Values of type $\Box A$ represent explicitly heap-allocated values of type A .

These primitives do not necessarily form an exhaustive list; indeed, we will sometimes expect further properties of these primitives to hold propositionally in the form of additional primitive lemmas, such as $(n : \text{Word}) \rightarrow \text{add } 0 \text{ } n =_{\text{Word}} n$. Precise details are only necessary when implementing such a system, and the definitions above are sufficient for the present discussion.

2.2 Extending λ_{PRIM} with inductive constructions

Building on top of λ_{PRIM} , we introduce an extension of the system for named inductive constructions in the meta-language, called λ_{IND} . First, we present a raw syntax for signatures Σ containing items Z consisting of data declarations, constructor declarations, and definitions:

$$\Sigma ::= \cdot \mid \Sigma, Z$$

$$Z ::= \text{data } \mathbf{D} \Delta : \uparrow \mathcal{U}_0 \mid \text{ctor } \mathbf{C} \Delta : \mathbf{D} \Delta \mid \text{closed } \mathbf{D} \vec{\mathbf{C}} \mid \text{def } \mathbf{f} : t = t$$

The symbols in blue represent labels, which are used to uniquely identify elements of a signature. A data definition $\text{data } \mathbf{D} \Delta : \uparrow \mathcal{U}_0$ introduces a new inductive type \mathbf{D} with a telescope Δ of arguments. Inductive types should be thought of as *codes* for object-level types, eventually to be replaced by the defined representations. A constructor definition $\text{ctor } \mathbf{C} \Delta_{\mathbf{C}} : \mathbf{D} \Delta$ introduces a new constructor \mathbf{C} for the inductive type \mathbf{D} , with a telescope $\Delta_{\mathbf{C}}$ of parameters which may depend on \mathbf{D} 's parameters. Closed declarations $\text{closed } \mathbf{D} \vec{\mathbf{C}}$ specify that the constructors $\vec{\mathbf{C}}$ are the only constructors for the inductive type \mathbf{D} in the present signature. A definition $\text{def } \mathbf{f} : t = t$ introduces a new symbol \mathbf{f} of type t and value t . The reason for including named definitions in a signature is to allow them to be overridden as part of data representations.

The system allows for the definition of inductive types as well as inductive families; the telescope Δ in a data declaration defines the parameters of the inductive type, and index refinement can occur in the constructor telescopes $\Delta_{\mathbf{C}}$ through the usage of equality types, paired with the ability to reference the parameters of the inductive type Δ . Such an approach is syntactically simpler than the standard approach of constructor signatures as Π types common in most proof assistants, but is equivalent in expressive power [1]. For example, to define the type of vectors indexed by

their length, we can write

```

data Vec ( $A : \mathbb{U}_0, n : \text{Nat}$ ) :  $\mathbb{U}_0$ 
ctor nil ( $n =_{\text{Nat}} z$ ) : Vec  $A$   $n$ 
ctor cons ( $a : A, m : \text{Nat}, v : \text{Vec } A \ m, m =_{\text{Nat}} s \ n$ ) : Vec  $A$   $n$ 
closed Vec (nil, cons)

```

In a real implementation, these definitions can be elaborated from a more familiar syntax involving index refinement.

The syntax of terms is extended accordingly with labels applied to arguments

$$t ::= \dots \mid \mathbf{L} \vec{t}$$

where \vec{t} denotes a sequence of terms and \mathbf{L} refers to any valid label in a signature (data, constructor, or definition). For example, given a constructor

```

ctor cons ( $x : T, xs : \text{List } T$ ) : List  $T$ 

```

we can write **cons** $a \ l$ to denote the full application of the ‘cons’ constructor to an element a and a list l . Partial applications are also allowed as syntactic syntactic sugar, so that **cons** a is shorthand for $l \mapsto \text{cons } a \ l$. Each inductive definition **data** $\mathbf{D} \Delta : \mathbb{U}_0$ exposes an additional label **case_D** which is used to perform case analysis on terms of the inductive type \mathbf{D} . As such, valid labelled application terms in a signature are given through the typing rules of λ_{IND} and do not correspond exactly to the labels *present* in the signature’s items.

In λ_{IND} , contexts are fibered over signatures, so that a context Γ is well formed with respect to a signature Σ . The syntax for telescopes Γ and contexts Δ is

$$\Gamma ::= \cdot \mid \Gamma, x : t$$

$$\Delta ::= \cdot \mid \Delta, x : t.$$

As is standard in 2LTT, contexts can be extended with types from any stage. Telescopes are very similar to contexts, but restricted to types from a single stage and well formed with respect to a context Γ , meaning that telescopes can contain open terms. We use the notation $\Delta \rightarrow t$ to denote a repeated function type with parameters from Δ and codomain t which may depend on the parameters. Additionally, we will sometimes explicitly bind the names of a telescope such as $(\vec{x} : \Delta) \rightarrow t[\vec{x}]$. Similar syntax is used to extend contexts with telescopes: Γ, Δ or $\Gamma, \vec{x} : \Delta$.

The language λ_{IND} is equipped with the following judgment forms:

- $\Sigma \text{ sig}$ – The signature Σ is well-formed.
- $\Sigma \vdash \Gamma \text{ con}$ – In signature Σ the context Γ is well-formed.
- $\Sigma \mid \Gamma \vdash \Delta \text{ tel}_i$ – In signature Σ and context Γ , the telescope Δ is well-formed in stage i .
- $\Sigma \mid \Gamma \vdash T \text{ type}_i$ – In signature Σ and context Γ , T is a well-formed type in stage i .
- $\Sigma \mid \Gamma \vdash a : T$ – In signature Σ and context Γ , a is a well-formed term of type T .
- $Z \in \Sigma$ – The item Z is present in the signature Σ .
- $\Sigma \vdash Z$ – The item Z is well-formed in the signature Σ .
- $\mathbf{L} \text{ label} \notin \Sigma$ – The label \mathbf{L} does not appear in the signature Σ .

First, we present the rules for well formed signatures, contexts and telescopes in fig. 1.

Next, the rules for well-formed items in signatures are given in fig. 2.

$$\begin{array}{c}
\text{SIG-EMPTY} \\
\frac{}{\cdot \text{sig}} \\
\text{SIG-EXTEND} \\
\frac{\Sigma \text{ sig} \quad \Sigma \vdash Z}{\Sigma, Z \text{ sig}} \\
\text{CON-EMPTY} \\
\frac{}{\Sigma \vdash \cdot \text{con}} \\
\text{CON-EXTEND} \\
\frac{\Sigma \vdash \Gamma \text{ con} \quad \Sigma \mid \Gamma \vdash T \text{ type}_i}{\Sigma \vdash \Gamma, T \text{ con}} \\
\text{TEL-EMPTY} \\
\frac{}{\Sigma \mid \Gamma \vdash \cdot \text{tel}_i} \\
\text{TEL-EXTEND} \\
\frac{\Sigma \mid \Gamma \vdash \Delta \text{ tel}_i \quad \Sigma \mid \Gamma \vdash T \text{ type}_i}{\Sigma \mid \Gamma \vdash \Gamma, T \text{ tel}_i}
\end{array}$$

Fig. 1. Rules for signatures, contexts and telescopes in λ_{IND} .

$$\begin{array}{c}
\text{DATA-ITEM} \\
\frac{\Sigma \mid \cdot \vdash \Delta \text{ tel}_1 \quad \text{D label} \notin \Sigma}{\Sigma \vdash \text{data } \text{D} \Delta : \uparrow \mathcal{U}_0} \\
\text{CTOR-ITEM} \\
\frac{\text{data } \text{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \text{closed } \text{D}, _ \notin \Sigma \quad \text{C label} \notin \Sigma \quad \Sigma \mid \Delta \vdash \Delta_{\text{C}} \text{ tel}_1}{\Sigma \vdash \text{ctor } \text{C} \Delta_{\text{C}} : \text{D} \Delta} \\
\text{CLOSED-ITEM} \\
\frac{\text{data } \text{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \forall i \in I. \text{ctor } \text{C}_i \Delta_{\text{C}_i} : \text{D} \Delta \in \Sigma \quad \text{closed } \text{D}, _ \notin \Sigma}{\Sigma \vdash \text{closed } \text{D}, \vec{\text{C}}} \\
\text{DEF-ITEM} \\
\frac{\Sigma \mid \cdot \vdash m : M \quad \text{f label} \notin \Sigma}{\Sigma \vdash \text{def } \text{f} : M = m}
\end{array}$$

Fig. 2. Rules for items in signatures in λ_{IND} .

$$\begin{array}{c}
\text{DATA-FORM} \\
\frac{\text{data } \text{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \Sigma \mid \Gamma \vdash \vec{t} : \Delta}{\Sigma \mid \Gamma \vdash \text{D} \vec{t} : \uparrow \mathcal{U}_0} \\
\text{DATA-INTRO} \\
\frac{\text{data } \text{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \text{ctor } \text{C} \Delta_{\text{C}} : \text{D} \Delta \in \Sigma \quad \Sigma \mid \Gamma \vdash \vec{t} : \Delta \quad \Sigma \mid \Gamma \vdash \vec{u} : \Delta_{\text{C}}[\vec{t}]}{\Sigma \mid \Gamma \vdash \text{C} \vec{u} : \text{D} \vec{t}} \\
\text{DATA-ELIM} \\
\frac{\text{data } \text{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \forall i \in I. \text{ctor } \text{C}_i \Delta_{\text{C}_i} : \text{D} \Delta \in \Sigma \quad \text{closed } \text{D} \vec{\text{C}} \in \Sigma \quad \Sigma \mid \Gamma \vdash \vec{t} : \Delta \quad \Sigma \mid \Gamma \vdash \eta : \text{D} \vec{t} \quad \Sigma \mid \Gamma, \vec{x} : \Delta, h : \text{D} \vec{x} \rightarrow T : \uparrow \mathcal{U}_0 \quad \forall i \in I. \Sigma \mid \Gamma, \vec{x} : \Delta, \vec{u} : \Delta_{\text{C}_i}[\vec{x}] \vdash m_i : T[\vec{x}, \text{C}_i \vec{u}]}{\Sigma \mid \Gamma \vdash \text{case}_{\text{D}} \eta \vec{m} : T[\vec{t}, \eta]} \\
\text{DEF-INTRO} \\
\frac{\text{def } \text{f} : M = m \in \Sigma}{\Sigma \mid \Gamma \vdash \text{f} : M}
\end{array}$$

Fig. 3. Typing rules for terms and types related to items in λ_{REP} .

There is a difference between an item being well-formed in a signature, and being *present* in a signature. The former describes items which are candidates to be added to a signature, which often relies on the absence of certain items in the signature, such as duplicate labels not being present, or a data type not being closed. We do not explicitly require that defined data types are well-founded or strictly positive, though this is a desirable property for a real implementation.

Actually, I think this should be a requirement, because otherwise case elimination is not quite valid...

Having defined well-formed signatures, it is now time to define the rules of well-formed types and terms in λ_{IND} , shown in fig. 3. For brevity, only the rules that relate to items in signatures are shown, that is, rules regarding the introduction and elimination of data types, constructors, and definitions. The rest of the rules are standard and can be found in the literature on 2LTT [3], with the modification that everything is now fibered over signatures. This completes a description of a staged language with dependent types and inductive constructions, λ_{IND} . In the next section, we move on to the main focus of this work, data representations.

2.3 Data representations in λ_{REP}

So far, we have described λ_{IND} , a dependently-typed staged language with object-level machine primitives as well as named inductive constructions and definitions. We now introduce data representations in the syntax for items, forming the language λ_{REP} . The goal of data representations is to provide a way to represent data types, constructions, and definitions in a more efficient manner, by transforming them into more suitable data structures for the target architecture. This is achieved by defining a new kind of item in signatures, called a representation, which specifies how to represent a given meta-level item in the object language. With no further restrictions, a user would be able to arbitrarily change the semantics of the meta-level items through representations, which would be undesirable. Instead, we restrict data representations to preserve the intended semantics of the original items, allowing for a correct-by-construction transformation from λ_{REP} to λ_{PRIM} . This is done in different ways for data types, constructors, and definitions.

First, we extend the raw syntax of items Z with representations, forming λ_{REP} :

$$Z ::= \dots \mid \mathbf{repr} \ L \ \Delta \ \mathbf{as} \ t$$

A representation $\mathbf{repr} \ L \ \Delta \ \mathbf{as} \ t$ asks for a definition L with parameters Δ to be represented by a term t . Data, constructor and definition representations all share the same raw syntax and are distinguished by the subject L . In the typing rules we expect that every meta-level type A has a defined representation $\mathcal{R}A$. Extensions to this system could support a sub-class of types with representations, retaining the ability to have purely meta-level types with no defined representations. We define $\mathcal{R}A$ to be the representation of a type A and $\mathcal{R}a$ the representation of a term a , whose type is $\mathcal{R}A$. These compute definitionally according to rules developed in section 3, and indeed constitute the elaboration process into λ_{PRIM} . For now, however, the definition of \mathcal{R} can remain opaque. With that, we can now define the additional typing rules for λ_{REP} , in fig. 4

The correctness of a representation of some data type D is ensured by the correctness of the representations of each of its constructors C_i and the case analysis function \mathbf{case}_D . Therefore there is no special correctness condition for the data representation type itself. The constructor representations $\mathbf{repr} \ C_i \ \mathcal{R}\Delta_{C_i} \ \mathbf{as} \ t_i$ form an *algebra* over the endofunctor $\llbracket D \rrbracket$ associated with the data type D , by packaging all the representation values $\pi_1 t_i$, where the carrier object is the defined representation type A . Moreover, this algebra is *injective* in the sense that the representation values are unique for each constructor, ensured by the equality constraints $\pi_1 t_i \neq \pi_1 t_j$ for $i \neq j$.

It is known that algebras over indexed inductive types can be interpreted as *ornaments* [2]; an inductive type is decorated with the values of the algebra at each node. We can apply the constructor representation algebra $(t_i \mid i \in I)$ to the inductive type D to obtain the ornamented type \tilde{D} . The case analysis \mathbf{case}_D representation is thus a *section* of \tilde{D} by the represented type A . In other words, it must ensure that the subject η of the case analysis is used to index into the ornamented type \tilde{D} , or put another way, that the propositional equality $\eta = \pi_1 t_i$ holds when the branch m_i is invoked. Overall, this yields an isomorphism between the inductive type D and the represented type A , in the Kleisli category of the code-generation monad of the object-level language (TODO: expand on this!).

$$\frac{\text{REPR-DEF} \quad \text{def } f : M = m \in \Sigma \quad \Sigma \mid \cdot \vdash a : (m' : \mathcal{RM}) \times (\mathcal{RM} =_{\mathcal{RM}} m')}{\Sigma \vdash \text{repr } f \text{ as } a}$$

Todo:

Extensions:

3 ELABORATION INTO A CORE STAGED LANGUAGE

- if **data** $D \Delta : \uparrow \mathcal{U}_0 \in \Sigma$, then $\exists A \vec{x}. \text{repr } D \vec{x} \text{ as } A \in \Sigma$,
- if **ctor** $C \Delta_C : D \Delta \in \Sigma$, then $\exists t \vec{z}. \text{repr } C \vec{z} \text{ as } t \in \Sigma$,
- if **closed** $D \vec{C} \in \Sigma$, then $\exists T \vec{a} \eta \vec{m}. \text{repr case}_D T \vec{a} \eta \vec{m} \text{ as } e \in \Sigma$, and
- if **def** $f : M = m \in \Sigma$, then $\exists a. \text{repr } f \text{ as } a \in \Sigma$.

Manuscript submitted to ACM

$$\begin{array}{lll}
\mathcal{R} : \text{Con } \Sigma \rightarrow \text{Con } \Sigma & \mathcal{R} : \text{Ty}_i \Sigma \Gamma \rightarrow \text{Ty}_i \Sigma \mathcal{R}\Gamma & \mathcal{R} : \text{Tm}_i \Sigma \Gamma T \rightarrow \text{Tm}_i \Sigma \mathcal{R}\Gamma \mathcal{R}T \\
\mathcal{R}(\cdot) = \cdot & \mathcal{R}(\mathbf{D} \vec{t}) = \uparrow \mathcal{R}A_{\mathbf{D}}[\mathcal{R}\vec{t}] & \mathcal{R}(\mathbf{C} \vec{t}) = \langle \mathcal{R}(\pi_1 t_{\mathbf{C}})[\mathcal{R}\vec{t}] \rangle \\
\mathcal{R}(\Gamma, T) = \mathcal{R}\Gamma, \mathcal{R}T & \text{else recurse with } \mathcal{R} & \mathcal{R}(\text{case}_{\mathbf{D}} \eta \vec{m}) = \langle \mathcal{R}ec[_{_}, _]\mathcal{R}\eta, \mathcal{R}\vec{m} \rangle \\
& & \mathcal{R}\mathbf{f} = \pi_1 a_{\mathbf{f}} \\
& & \text{else recurse with } \mathcal{R}
\end{array}$$

Fig. 5. Definition of \mathcal{R} , where Σ is a concrete signature.

4 RECOVERING THE DEFAULT LINKED REPRESENTATION

DEFAULT-DATA

Σ concrete

$$\begin{array}{c}
\Sigma \vdash Z := \mathbf{data} \mathbf{D} \Delta : \uparrow \mathcal{U}_0 \quad \forall i \in I. \Sigma, Z, (X_j)_{j < i} \vdash X_i := \mathbf{ctor} \mathbf{C}_i \Delta_{\mathbf{C}_i} : \mathbf{D} \Delta \quad \Sigma, Z, (X_i)_{i \in I} \vdash L := \mathbf{closed} \mathbf{D} \vec{\mathbf{C}} \\
\hline
\Sigma, Z, X, L \vdash \mathbf{repr} \mathbf{D} \vec{x} \mathbf{as} (\mu F. \vec{x} \mapsto (c : \text{Word}_{\llbracket I \rrbracket}) \times \square(\text{switch } c (\text{Ctors } \vec{x} F)))) \vec{x} \\
\text{where } \text{Ctors } \vec{x} F = ((d : \text{Data}_i \vec{x}) \times ((r : \text{Rec}_i \vec{x} d) \rightarrow F(\text{Idx}_i r)))_{i \in I}
\end{array}$$

REPR-CTOR

$$\begin{array}{c}
\mathbf{data} \mathbf{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \forall i \in I. \mathbf{ctor} \mathbf{C}_i \Delta_{\mathbf{C}_i} : \mathbf{D} \Delta \in \Sigma \quad \mathbf{repr} \mathbf{D} \mathcal{R} \Delta \mathbf{as} A \in \Sigma \\
\forall j < i. \mathbf{repr} \mathbf{C}_j \mathcal{R} \Delta_{\mathbf{C}_j} \mathbf{as} t_j \in \Sigma \quad \Sigma \mid \vec{x} : \mathcal{R} \Delta, \vec{z} : \mathcal{R} \Delta_{\mathbf{C}_i} \vdash t_i : (a : A) \times \prod_{j < i} ((\vec{y} : \mathcal{R} \Delta_{\mathbf{C}_j}) \rightarrow a \neq \pi_1 t_j[\vec{x}, \vec{y}]) \\
\hline
\Sigma \vdash \mathbf{repr} \mathbf{C}_i \vec{z} \mathbf{as} t_i
\end{array}$$

REPR-CASE

$$\begin{array}{c}
\mathbf{data} \mathbf{D} \Delta : \uparrow \mathcal{U}_0 \in \Sigma \quad \mathbf{repr} \mathbf{D} \mathcal{R} \Delta \mathbf{as} A \in \Sigma \quad \forall i \in I. \mathbf{ctor} \mathbf{C}_i \Delta_{\mathbf{C}_i} : \mathbf{D} \Delta \in \Sigma \quad \forall i \in I. \mathbf{repr} \mathbf{C}_i \mathcal{R} \Delta_{\mathbf{C}_i} \mathbf{as} t_i \in \Sigma \\
\mathbf{closed} \mathbf{D} \vec{\mathbf{C}} \in \Sigma \quad \Sigma \mid T : (\vec{x} : \mathcal{R} \Delta) \rightarrow \uparrow A[\vec{x}] \rightarrow \mathcal{U}_0, \vec{a} : \mathcal{R} \Delta, \eta : A[\vec{a}], \vec{m} : \text{Cases } T \vec{a} \eta \vdash e : T(\vec{a}, \eta) \\
\text{where } \text{Cases } T \vec{a} \eta = (m_i : (\vec{y} : \mathcal{R} \Delta_{\mathbf{C}_i}) \rightarrow (\eta =_{A[\vec{a}]} \pi_1 t_i[\vec{a}, \vec{y}]) \rightarrow T(\vec{a}, \pi_1 t_i[\vec{a}, \vec{y}]) \mid i \in I) \\
\hline
\Sigma \vdash \mathbf{repr} \text{case}_{\mathbf{D}} T \vec{a} \eta \vec{m} \mathbf{as} e
\end{array}$$

DEFAULT-DEF

$$\begin{array}{c}
\Sigma \text{ concrete} \quad \Sigma \vdash F := \mathbf{def} \mathbf{f} : M = m \\
\hline
\Sigma, F \vdash \mathbf{repr} \mathbf{f} \mathbf{as} (\mathcal{R}m, \mathbf{refl})
\end{array}$$

Fig. 6. Default data representations in λ_{rep} . These are not rules, but derivations.

5 NATURAL NUMBERS, AND OTHER EXAMPLES

6 CONCLUSIONS AND FUTURE WORK

REFERENCES

- [1] Jesper Cockx and Andreas Abel. 2018. Elaborating dependent (co)pattern matching. *Proc. ACM Program. Lang.* 2, ICFP (July 2018), 1–30. <https://doi.org/10.1145/3236770>
- [2] Pierre-Evariste Dagand. 2017. The essence of ornaments. *J. Funct. Programming* 27 (Jan. 2017), e9. <https://www.cambridge.org/core/services/aop-cambridge-core/content/view/4D2DF6F4FE23599C8C1FEA6C921A3748/S0956796816000356a.pdf/div-class-title-the-essence-of-ornaments-div.pdf>
- [3] András Kovács. 2022. Staged compilation with two-level type theory. (Sept. 2022). arXiv:2209.09729 [cs.PL]