

# Program Term Generation Through Enumeration of Indexed datatypes (Thesis Proposal)

Cas van der Rest

January 24, 2019

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Problem Statement . . . . .	2
1.2	Research Questions and Contributions . . . . .	2
<b>2</b>	<b>Background</b>	<b>2</b>
2.1	Dependently Typed Programming & Agda . . . . .	3
2.1.1	Dependent Type Theory . . . . .	3
2.1.2	Curry-Howard Isomorphism . . . . .	3
2.1.3	Codata . . . . .	3
2.2	Property Based Testing . . . . .	3
2.2.1	Existing Libraries . . . . .	3
2.2.2	Generating Test Data . . . . .	3
2.3	Generic Programming & Type Universes . . . . .	3
2.3.1	Regular Datatypes . . . . .	4
2.3.2	Ornaments . . . . .	5
2.3.3	Functorial Species . . . . .	7
2.3.4	Indexed Functors . . . . .	7
2.4	Blockchain Semantics . . . . .	7
2.4.1	BitML . . . . .	7
2.4.2	UTXO & Extended UTXO . . . . .	7
<b>3</b>	<b>Preliminary results</b>	<b>8</b>
<b>4</b>	<b>Timetable and planning</b>	<b>8</b>

## 1 Introduction

A common way of asserting a program's correctness is by defining properties that should universally hold, and asserting these properties over a range of random inputs. This technique is commonly referred to as *property based testing*, and generally consists of a two-step process. Defining properties that universally hold on all inputs, and defining *generators* that sample random values from the space of possible inputs. *QuickCheck* [3] is likely the most well known tool for performing property based tests on haskell programs.

Although coming up with a set of properties that properly captures a program's behaviour might initially seem to be the most involved part of the process, defining suitable generators for complex input data is actually quite difficult as well. Questions such as how to handle datatypes that are inhabited by an infinite number of values arise, or how to deal with constrained input data. The answers to these questions are reasonably well understood for *Algebraic datatypes (ADT's)*, but no general solution exists when more complex input data is required. In particular, little is known about enumerating and generating inhabitants of *Indexed datatypes*.

The latter may be of interest when considering property based testing in the context of languages with a more elaborate type system than Haskell's, such as *Agda* or *Idris*. Since the techniques used in existing tools such as QuickCheck and SmallCheck for the most part only apply to regular datatypes, meaning that there is no canonical way of generating inhabitants for a large class of datatypes in these languages.

Besides the obvious applications to property based testing in the context of dependently typed languages, a broader understanding of how we can generate inhabitants of indexed datatypes may prove useful in other areas as well. Since we can often capture a programming language's semantics as an indexed datatype, efficient generation of inhabitants of such a datatype may prove useful for testing compiler infrastructure.

### 1.1 Problem Statement

### 1.2 Research Questions and Contributions

What is the problem? Illustrate with an example. [1, 12]

What is/are your research questions/contributions? [3]

## 2 Background

What is the existing technology and literature that I'll be studying/using in my research [6, 10, 11, 14]

## 2.1 Dependently Typed Programming & Agda

### 2.1.1 Dependent Type Theory

Dependent type theory extends a type theory with the possibility of defining types that depend on values. In addition to familiar constructs, such as the unit type ( $\top$ ) and the empty type  $\perp$ , one can use so-called  $\Pi$ -types and  $\Sigma$ -types.  $\Pi$ -types capture the idea of dependent function types, that is, *functions* whose output type may depend on the values of its input. Given some type  $A$  and a family  $P$  of types indexed by values of type  $A$  (i.e.  $P$  has type  $A \rightarrow \text{Type}$ ),  $\Pi$ -types have the following definition:

$$\Pi_{(x:A)} P(x) \quad \equiv \quad (x : A) \rightarrow P(x)$$

In a similar spirit,  $\Sigma$ -types are ordered *pairs* of which the type of the second value may depend on the first value of the pair.

$$\Sigma_{(x:A)} P(x) \quad \equiv \quad (x : A) \times P(x)$$

The Curry-Howard equivalence extends to  $\Pi$ - and  $\Sigma$ -types as well: they can be used to model universal and existential quantification.

### 2.1.2 Codata

Agda requires all functions to be total. This means that they are required to yield a result on all inputs in a *finite* amount of time. This means that we cannot work with infinite structures in the same way as in Haskell. For example, the following is perfectly fine in Haskell:

```
infinity :: Nat
infinity = Suc infinity

isSmaller :: Nat → Nat → Bool
isSmaller (Suc n) Zero    = False
isSmaller Zero   (Suc n) = True
isSmaller (Suc n) (Suc m) = isSmaller n m
```

## 2.2 Property Based Testing

### 2.2.1 Existing Libraries

### 2.2.2 Generating Test Data

## 2.3 Generic Programming & Type Universes

If we desire to abstract over the structure of datatypes, we need a suitable type universe to do so. Many such universes have been developed and studied; this section discusses a few of them.

### 2.3.1 Regular Datatypes

The term *regular datatypes* is often used to refer to the class of datatypes that can be assembled using any combination of products, coproducts, unary constructors, constants (a position that is inhabited by a value of another type) and recursive positions. Roughly, this class consists of ADT's in haskell, though mutual recursion is not accounted for.

Any value that lives in the induced by these combinators describes a regular datatype, and is generally referred to as a *pattern functor*. We can define a datatype in agda that captures these values:

```

data Reg : Set → Set where
  U      : Reg ⊥
  K      : (a : Set) → Reg a
  _⊕_    : ∀ {a : Set} → Reg a → Reg a → Reg a
  _⊗_    : ∀ {a : Set} → Reg a → Reg a → Reg a
  I      : Reg ⊥

```

Pattern functors, i.e. values of the *Reg* datatype, can be interpreted as types. Inhabitants of the interpretation of a pattern functor correspond to the inhabitants of the type that is represented by said pattern functor. We use the following interpretation function:

```

[[_]] : Reg → Set → Set
[[ U      ]] r = ⊤
[[ K a    ]] r = a
[[ reg1 ⊕ reg2 ]] r = [[ reg1 ]] r ⊔ [[ reg2 ]] r
[[ reg1 ⊗ reg2 ]] r = [[ reg1 ]] r × [[ reg2 ]] r
[[ I      ]] r = r

```

Notice that recursive positions are left explicit. This means that we require an appropriate fixed-point combinator to find a pattern functor's representation in **Set**.

```

data μ (f : Reg) : Set where
  'μ : [[ f ]] (μ f) → μ f

```

**Example** Consider the pattern functor corresponding to the definition of *List*:

```

List' : Set → Set
List' a = μ (U ⊕ (K a ⊗ I))

```

Notice that this pattern functor denotes a choice between a unary constructor (`[]`), and a constructor that takes a constant of type *a* and a recursive positions as arguments (`::`). We can define conversion functions between the standard *List* type, and the interpretation of our pattern functor:

```

fromList : ∀ {a : Set} → List a → List' a
fromList [] = 'μ (inj1 tt)
fromList (x :: xs) = 'μ (inj2 (x , fromList xs))

```

```

toList : ∀ {a : Set} → List' a → List a
toList ('μ (inj1 tt)) = []
toList ('μ (inj2 (fst , snd))) = fst :: toList snd

```

With these definitions, it is now trivial to show that there is indeed an isomorphism between the two:

```

isoList1 : ∀ {a : Set} {xs : List a} → toList (fromList xs) ≡ xs
isoList1 {xs = []} = refl
isoList1 {xs = x :: xs} = cong (_ :: _ x) isoList1

```

```

isoList2 : ∀ {a : Set} {xs : List' a} → fromList (toList xs) ≡ xs
isoList2 {xs = 'μ (inj1 x)} = refl
isoList2 {xs = 'μ (inj2 (fst , snd))} = cong ('μ ∘ inj2 ∘ _, _ fst) isoList2

```

Using such isomorphisms, we can automatically derive functionality for datatypes that can be captured using pattern functors. We will see an example of this in section 3, where we will derive enumeration of inhabitants for arbitrary pattern functors.

### 2.3.2 Ornaments

*Ornaments* [5] provide a type universe in which we can describe the structure of indexed datatypes in a very index-centric way. Indexed datatypes are described by *Signatures*, consisting of three elements:

- A function  $Op : I \rightarrow Set$ , that relates indices to operations/constructors
- A function  $Ar : Op\ i \rightarrow Set$ , that describes the arity (with respect to recursive positions) for an operation
- A typing discipline  $Ty : Ar\ op \rightarrow I$ , that describes indices for recursive positions.

When combined into a single structure, we say that  $\Sigma_D$  gives the signature of some indexed datatype  $D : I \rightarrow Set$ :

$$\Sigma_D(I) = \begin{cases} Op : I \rightarrow Set \\ Ar : Op\ i \rightarrow Set \\ Ty : Ar\ op \rightarrow I \end{cases}$$

**Example** Let us consider the signature for the *Vec* type, given by  $\Sigma_{Vec}(\mathbb{N})$ . Recall the definition of the *Vec* datatype in listing 1. It has the following relation between index and operations:

$$\begin{aligned} \text{Op-vec} &: \forall \{a : \text{Set}\} \rightarrow \mathbb{N} \rightarrow \text{Set} \\ \text{Op-vec zero} &= \top \\ \text{Op-vec } \{a\} (\text{suc } n) &= a \end{aligned}$$

If the index is *zero*, we have only the unary constructor  $[]$  at our disposal, hence  $\text{Op-vec zero} = \text{top}$ . If the index is *suc n*, the number of possible constructions for *Vec* corresponds to the set of inhabitants of its element type, hence we say that  $\text{Op-vec (suc } n) = a$ .

The  $[]$  constructor has no recursive argument, so its arity is  $\perp$ . Similarly, *cons a* takes one recursive argument, so its arity is  $\top$ :

$$\begin{aligned} \text{Ar-vec} &: \forall \{a : \text{Set}\} \rightarrow (n : \mathbb{N}) \rightarrow \text{Op-vec } \{a\} n \rightarrow \text{Set} \\ \text{Ar-vec zero tt} &= \perp \\ \text{Ar-vec (suc } n) \text{ op} &= \top \end{aligned}$$

The definition of  $::$  dictates that if the index is equal to *suc n*, the index of the recursive argument needs to be *n*. We interpret this as follows: if a vector has length (suc n), its tail has length n. This induces the following typing discipline for *Vec*:

$$\begin{aligned} \text{Ty-vec} &: \forall \{a : \text{Set}\} \rightarrow (n : \mathbb{N}) \rightarrow (\text{op} : \text{Op-vec } \{a\} n) \rightarrow \text{Ar-vec } n \text{ op} \rightarrow \mathbb{N} \\ \text{Ty-vec zero } a \text{ } () & \\ \text{Ty-vec (suc } n) \text{ a tt} &= n \end{aligned}$$

This defines the signature for *Vec*:  $\Sigma_{Vec} \triangleq \text{Op-vec} \triangleleft^{\text{Ty-vec}} \text{Ar-vec}$ .

---

```

data Vec {a} (A : Set a) :  $\mathbb{N} \rightarrow \text{Set a}$  where
  [] : Vec A zero
  _::_ :  $\forall \{n\} (x : A) (xs : \text{Vec A } n) \rightarrow \text{Vec A (suc } n)$ 

```

---

**Listing 1:** Definition of *Vec*

We can define signatures for non-indexed datatypes as well by choosing a trivial index, e.g.  $I = \top$ . This gives  $\Sigma_{\mathbb{N}} \triangleq \text{Op-nat} \triangleleft^{\text{Ty-nat}} \text{Ar-nat}$  with the definitions given in listing 2:

---


$$\begin{aligned}\text{Op-nat} &: \top \rightarrow \text{Set} \\ \text{Op-nat tt} &= \top \uplus \top\end{aligned}$$

$$\begin{aligned}\text{Ar-nat} &: \text{Op-nat tt} \rightarrow \text{Set} \\ \text{Ar-nat} (\text{inj}_1 x) &= \perp \\ \text{Ar-nat} (\text{inj}_2 y) &= \top\end{aligned}$$

$$\begin{aligned}\text{Ty-nat} &: (\text{op} : \text{Op-nat tt}) \rightarrow \text{Ar-nat op} \rightarrow \top \\ \text{Ty-nat} (\text{inj}_1 x) () & \\ \text{Ty-nat} (\text{inj}_2 y) \text{tt} &= \text{tt}\end{aligned}$$


---

**Listing 2:** Signature definition for  $\mathbb{N}$

### 2.3.3 Functorial Species

### 2.3.4 Indexed Functors

## 2.4 Blockchain Semantics

### 2.4.1 BitML

### 2.4.2 UTXO & Extended UTXO

- Libraries for property based testing (QuickCheck, (Lazy) SmallCheck, QuickChick, QuickSpec)
- Type universes (ADT's, Ornaments) [5, 8]
- Generic programming techniques. (pattern functors, indexed functors, functorial species)
- Techniques to generate complex or constrained data (Generating constrained random data with uniform distribution, Generators for inductive relations)
- Techniques to speed up generation of data (Memoization, FEAT)
- Formal specification of blockchain (bitml, (extended) UTxO ledger) [15, 16]
- Representing potentially infinite data in Agda (Colists, coinduction, sized types)

Below is a bit of Agda code:

---

```

 $\Gamma\text{-match} : (\tau : \text{Ty}) \rightarrow \langle\langle \omega_i (\lambda \Gamma \rightarrow \Sigma[\alpha \in \text{Id}] \Gamma [\alpha \mapsto \tau]) \rangle\rangle$ 
 $\Gamma\text{-match } \tau \mu \emptyset = \text{uninhabited}$ 
 $\Gamma\text{-match } \tau \mu (\alpha \mapsto \sigma :: \Gamma) \text{ with } \tau \stackrel{?}{=} \sigma$ 
 $\Gamma\text{-match } \tau \mu (\alpha \mapsto \tau :: \Gamma) \mid \text{yes refl} = \llbracket (\alpha, \text{TOP}) \rrbracket$ 
 $\qquad \qquad \qquad \parallel \llbracket (\Sigma\text{-map POP}) (\mu \Gamma) \rrbracket$ 
 $\Gamma\text{-match } \tau \mu (\alpha \mapsto \sigma :: \Gamma) \mid \text{no } \neg p = \llbracket (\Sigma\text{-map POP}) (\mu \Gamma) \rrbracket$ 

```

---

**Listing 3:** Definition of  $\Gamma\text{-match}$ 


---

```

data Env : Set where
   $\emptyset : \text{Env}$ 
   $\_ \mapsto \_ :: \_ : \text{Id} \rightarrow \text{Ty} \rightarrow \text{Env} \rightarrow \text{Env}$ 

data  $\_[_ \mapsto \_]$  : Env  $\rightarrow$  Id  $\rightarrow$  Ty  $\rightarrow$  Set where

  TOP :  $\forall \{ \Gamma \alpha \tau \}$ 
         $\rightarrow (\alpha \mapsto \tau :: \Gamma) [\alpha \mapsto \tau]$ 

  POP :  $\forall \{ \Gamma \alpha \beta \tau \sigma \}$   $\rightarrow \Gamma [\alpha \mapsto \tau]$ 
         $\rightarrow (\beta \mapsto \sigma :: \Gamma) [\alpha \mapsto \tau]$ 

```

---

**Listing 4:** Environment definition and membership in *Agda*

### 3 Preliminary results

What examples can you handle already? [9]

What prototype have I built? [4, 7]

How can I generalize these results? What problems have I identified or do I expect? [13]

### 4 Timetable and planning

What will I do with the remainder of my thesis? [2]

Give an approximate estimation/timetable for what you will do and when you will be done.



---


$$\begin{array}{c}
TOP \frac{}{(a \mapsto t : \Gamma)[a \mapsto t]} \qquad POP \frac{\Gamma[a \mapsto t]}{(b \mapsto s : \Gamma)[a \mapsto t]} \\
\\
VAR \frac{\Gamma[a \mapsto \tau]}{\Gamma \vdash a : \tau} \qquad ABS \frac{\Gamma, a \mapsto \sigma \vdash t : \tau}{\Gamma \vdash \lambda a \rightarrow t : \sigma \rightarrow \tau} \\
\\
APP \frac{\Gamma \vdash f : \sigma \rightarrow \tau \quad \Gamma \vdash x : \sigma}{\Gamma \vdash fx : \tau} \qquad LET \frac{\Gamma \vdash e : \sigma \quad \Gamma, a \mapsto \sigma \vdash t : \tau}{\Gamma \vdash \text{let } a := e \text{ in } t : \tau}
\end{array}$$


---

**Listing 5:** Semantics of the *Simply Typed Lambda Calculus*

## References

- [1] Thorsten Altenkirch and Conor McBride. Generic programming within dependently typed programming. In *Generic Programming*, pages 1–20. Springer, 2003.
- [2] Koen Claessen, Jonas Duregård, and Michał H Pałka. Generating constrained random data with uniform distribution. *Journal of functional programming*, 25, 2015.
- [3] Koen Claessen and John Hughes. Quickcheck: a lightweight tool for random testing of haskell programs. *Acm sigplan notices*, 46(4):53–64, 2011.
- [4] Koen Claessen, Nicholas Smallbone, and John Hughes. Quickspec: Guessing formal specifications using testing. In *International Conference on Tests and Proofs*, pages 6–21. Springer, 2010.
- [5] Pierre-Évariste Dagand. The essence of ornaments. *Journal of Functional Programming*, 27, 2017.
- [6] Maxime Dénès, Catalin Hritcu, Leonidas Lampropoulos, Zoe Paraskevopoulou, and Benjamin C Pierce. Quickchick: Property-based testing for coq. In *The Coq Workshop*, 2014.
- [7] Jonas Duregård, Patrik Jansson, and Meng Wang. Feat: functional enumeration of algebraic types. *ACM SIGPLAN Notices*, 47(12):61–72, 2013.
- [8] Hsiang-Shang Ko and Jeremy Gibbons. Programming with ornaments. *Journal of Functional Programming*, 27, 2016.
- [9] Leonidas Lampropoulos, Zoe Paraskevopoulou, and Benjamin C Pierce. Generating good generators for inductive relations. *Proceedings of the ACM on Programming Languages*, 2(POPL):45, 2017.

- [10] Andres Löb and José Pedro Magalhaes. Generic programming with indexed functors. In *Proceedings of the seventh ACM SIGPLAN workshop on Generic programming*, pages 1–12. ACM, 2011.
- [11] Ulf Norell. Dependently typed programming in agda. In *International School on Advanced Functional Programming*, pages 230–266. Springer, 2008.
- [12] Colin Runciman, Matthew Naylor, and Fredrik Lindblad. Smallcheck and lazy smallcheck: automatic exhaustive testing for small values. In *Acm sigplan notices*, volume 44, pages 37–48. ACM, 2008.
- [13] Alexey Rodriguez Yakushev, Stefan Holdermans, Andres Löb, and Johan Jeuring. Generic programming with fixed points for mutually recursive datatypes. In *ACM Sigplan Notices*, volume 44, pages 233–244. ACM, 2009.
- [14] Brent A Yorgey. Species and functors and types, oh my! In *ACM Sigplan Notices*, volume 45, pages 147–158. ACM, 2010.
- [15] Joachim Zahnentferner. An abstract model of utxo-based cryptocurrencies with scripts. *IACR Cryptology ePrint Archive*, 2018:469, 2018.
- [16] Joachim Zahnentferner and Input Output HK. Chimeric ledgers: Translating and unifying utxo-based and account-based cryptocurrencies. Technical report, Cryptology ePrint Archive, Report 2018/262, 2018. <https://eprint.iacr.org/> . . . , 2018.