

Generic Enumerators

Cas van der Rest
c.r.vanderrest@students.uu.nl
Universiteit Utrecht

Wouter Swierstra
w.s.swierstra@uu.nl
Universiteit Utrecht

Manuel Chakravarty
manuel.chakravarty@iohk.io
...

Introduction

Since the introduction of QuickCheck [3], *property based testing* has proven to be effective for the discovery of bugs. However, defining the properties to test is only part of the story: it is equally important to *generate* suitable test data. In particular, requiring random test data to satisfy arbitrary preconditions can lead to skewed distributions: for example, naively generating random sorted lists will rarely yield long lists. As a result, developers need to design custom generators carefully—but these generators can become arbitrarily complex. When testing a compiler, for example, it can be quite challenging to define a good generator that is guaranteed to produce well-formed programs. [2, 6]

In this brief abstract we propose to address this problem using the observation that well-formed inputs can often be described by (indexed) inductive datatypes. By defining a *generic* procedure for *enumerating* indexed datatypes, we can obtain a way of safely generating precise test data.

Defining generators We will sketch how to define a generic enumerator for a collection of datatypes in several steps:

- We begin by defining some universe of types U together with its semantics of the form $\llbracket _ \rrbracket : U \rightarrow \text{Set}$;
- Next, we define a datatype generic function

$$\text{enumerate} : (u : U) \rightarrow (n : \mathbb{N}) \rightarrow \text{List } \llbracket u \rrbracket$$

This function produces a list of elements, bounded by some size parameter n ;

- Finally, we formulate the key *completeness* property that we expect of our enumerators:

$$\begin{aligned} \forall u \rightarrow (x : \llbracket u \rrbracket) \rightarrow \\ \Sigma [n \in \mathbb{N}] (x \in \text{enumerate } u \ n) \end{aligned}$$

This property states that for each possible x , there is some size n such that x occurs in our enumeration.

We will now sketch three increasingly complex universes, together with their corresponding generic enumerations.

Enumeration of regular types

One of the simplest universes that describes a wide class of algebraic datatypes is the *universe of regular types*. This universe contains the unit type, empty type, constant types, and is closed under products and coproducts.

data $\text{Reg} : \text{Set}$ **where**
 $\text{Z U I} : \text{Reg}$
 $_ \oplus _ \otimes _ : \text{Reg} \rightarrow \text{Reg} \rightarrow \text{Reg}$
 $\text{K} : \text{Set} \rightarrow \text{Reg}$

The associated semantics, $\llbracket _ \rrbracket : \text{Reg} \rightarrow \text{Set} \rightarrow \text{Set}$, maps values of type Reg to their corresponding pattern functor. By taking the fixpoint of such a pattern functor, we have a uniform representation of a wide class of (recursive) algebraic datatypes:

data $\text{Fix } (c : \text{Reg}) : \text{Set}$ **where**
 $\text{In} : \llbracket c \rrbracket (\text{Fix } c) \rightarrow \text{Fix } c$

Examples of regular types and their respective codes include natural numbers ($\text{U} \oplus \text{I}$) or lists ($\text{U} \oplus (\text{K } a \otimes \text{I})$).

It is reasonably straightforward to define a generic enumeration function:

$$\text{enumerate} : (c : \text{Reg}) \rightarrow (n : \mathbb{N}) \rightarrow \text{List } \llbracket \text{Fix } c \rrbracket$$

For example, the enumeration of a coproduct of two codes is a fair merge of the two recursive calls, and for products we enumerate all possible combinations of values.

Enumeration of Indexed Containers

What happens when we consider *indexed* datatypes? Initially, we will consider *indexed containers* [1, 5], a subset of all possible indexed types that are defined by induction over the index type.

Following the presentation in [5], we define indexed containers as a triple of *operations*, *arities* and *typing*:

$\text{Op} : i \rightarrow \text{Reg}$
 $\text{Ar} : \forall \{x\} \rightarrow \text{Fix } (\text{Op } x) \rightarrow \text{Reg}$
 $\text{Ty} : \forall \{x\} \{ \text{op} : \text{Fix } (\text{Op } x) \} \rightarrow \text{Fix } (\text{Ar } \text{op}) \rightarrow i$

The set $\text{Op } i$ describes the set of available operations at index i ; $\text{Ar } \text{op}$ the arities of each constructor; finally, $\text{Ty } \text{ar}$ gives the index corresponding to the recursive subtree at arity ar . Signatures are interpreted as a function from index to dependent pair, with the first element of the pair denoting a choice of constructor, and the second element being a function that maps each recursive subtree to a value of the type that results from applying the recursive argument with the index given by the typing discipline for that arity.

$$\begin{aligned} \llbracket \text{Op } \triangleleft \text{Ar} \mid \text{Ty} \rrbracket x &= \lambda i \rightarrow \\ &\Sigma [\text{op} \in \text{Fix } (\text{Op } i)] (\text{ar} : \text{Fix } (\text{Ar } \text{op})) \rightarrow x (\text{Ty } \text{ar}) \end{aligned}$$

Interpretations of signatures live in $I \rightarrow \text{Set}$, hence we also need adapt our fixpoint, Fix , accordingly.

Examples Many familiar indexed datatypes can be described using the universe of indexed containers, such as finite types (Fin), well-scoped lambda terms, or the type of vectors given below:

$\Sigma\text{-vec } a =$

```

let op-vec = ( $\lambda \{ \text{zero} \rightarrow U; (\text{suc } n) \rightarrow K \ a \}$ )
      ar-vec = ( $\lambda \{ \{ \text{zero} \} \text{tt} \rightarrow Z; \{ \text{suc } n \} \ x \rightarrow U \}$ )
      ty-vec = ( $\lambda \{ \{ \text{suc } n \} \{ a \} (\text{In tt}) \rightarrow n \}$ )
in op-vec  $\triangleleft$  ar-vec | ty-vec

```

Each index is associated with a unique operation. We map $\text{suc } n$ to a constant type in op-vec , since the $::$ constructor stores a value along its recursive subtree. The empty vector, $[]$, has no recursive subtrees, hence its arity is the empty type. Any non-empty vector has one subtree, so we assign its arity to be the unit type. This single subtree has an index that is one less than the original index, as described by ty-vec .

Generic enumerators In the definition of indexed containers, we carefully restricted the type of operations and arities to the universe of regular types. As a result, we can reuse the enumeration of regular types to write a generic enumerator for indexed containers. The second component of a signature's interpretation is a function type, so we require an enumerator for function types. Inspired by SmallCheck [7] we can define such an enumerator:

```

cogenerate :
  ( $\mathbb{N} \rightarrow \text{List } a$ )  $\rightarrow$  ( $c : \text{Reg}$ )  $\rightarrow \mathbb{N} \rightarrow \text{List } ([\![c]\!] \rightarrow a)$ 

```

This allows us to define enumerators for both components of the dependent pair:

```

enumOp :  $\forall i \rightarrow \mathbb{N} \rightarrow \text{List } ([\![\text{Op } i]\!])$ 
enumAr :  $\forall i \rightarrow (x : [\![\text{Op } i]\!]) \rightarrow \mathbb{N} \rightarrow \text{List } ([\![y : \text{Ar } x] \rightarrow r (\text{Ty } y)] \rightarrow r)$ 

```

We then sequence these operations using the monadic structure of lists:

```

 $\lambda n \rightarrow \text{enumOp } n \gg (\lambda \text{op} \rightarrow \text{op}, \text{enumAr } n \text{ op})$ 

```

Intuitively, this defines the enumeration of a signature as the union of the enumerations of its constructors.

Indexed Descriptions

Not all indexed families may be readily described as indexed containers. Consider, for example, the type of binary trees indexed by their number of nodes:

```

data Tree ( $a : \text{Set}$ ) :  $\mathbb{N} \rightarrow \text{Set}$  where
  Leaf : Tree  $a$  0
  Node :  $\forall \{n\ m\} \rightarrow \text{Tree } a \ n \rightarrow a \rightarrow \text{Tree } a \ m$ 
         $\rightarrow \text{Tree } a \ (\text{suc } (n + m))$ 

```

Without introducing further equalities, it is hard to capture the decomposition of the index $\text{suc } (n + m)$ into two subtrees of size n and m .

The universe of indexed descriptions, as defined in [4], is capable of representing arbitrary indexed families. This makes two key modifications to the universe of regular types:

firstly, recursive positions must store an additional field corresponding to their index. Secondly, a new combinator, Σ is added.

```

I : ( $i : I$ )  $\rightarrow \text{IDesc } I$ 
 $\Sigma : (S : \text{Set}) \rightarrow (T : S \rightarrow \text{IDesc } I) \rightarrow \text{IDesc } I$ 

```

Their interpretation is rather straightforward.

```

 $[\![I\ i]\!] r = r\ i$ 
 $[\![\Sigma\ S\ T]\!] r = \Sigma [s \in S] [\![T\ s]\!] r$ 

```

With the added Σ and var , we can now describe the Tree datatype:

```

tree :  $\text{Set} \rightarrow \mathbb{N} \rightarrow \text{IDesc } \mathbb{N}$ 
tree a zero = '1
tree a ( $\text{suc } n$ ) = ' $\Sigma (\Sigma (\mathbb{N} \times \mathbb{N}) \lambda \{(n', m') \rightarrow n' + m' \equiv n\})$ 
   $\lambda \{((n, m), \text{refl}) \rightarrow I\ n \otimes K\ a \otimes I\ m\}$ 

```

The dependency between the indices of the left- and right subtrees of nodes is captured by having their description depend on a pair of natural numbers together with a proof that they sum to the required index.

Enumerators for indexed descriptions Since the IDesc universe largely exposes the same combinators as the universe of regular types, we only really need to define enumerate for the Σ combinator. This is straightforward once we can enumerate its first component.

```

enumerate : ( $\delta : \text{IDesc } I$ )  $\rightarrow \mathbb{N} \rightarrow \text{List } ([\![\delta]\!])$ 
enumerate ( $\Sigma\ s\ g$ ) =
   $\lambda n \rightarrow \text{gen } n \gg (\lambda x \rightarrow x, \text{enumerate } (g\ s)\ n)$ 

```

However, since Σ may range over any type in Set , we have no generic procedure to obtain a suitable enumerator. This creates a separation between the parts of a datatype for which an enumerator can be assembled mechanically, and those parts for which this would be too difficult.

In the case of the Tree datatype, we see that those elements that make it hard to generically enumerate inhabitants of this datatype emerge quite naturally; we merely need to supply a generator inverting addition:

```

 $+^{-1} : (n : \mathbb{N}) \rightarrow \mathbb{N}$ 
 $\rightarrow \text{List } (\Sigma (\mathbb{N} \times \mathbb{N}) \lambda \{(n', m') \rightarrow n' + m' \equiv n\})$ 

```

enumerate needs nothing beyond $+^{-1}$ in order to be able to enumerate inhabitants of Tree.

TODO: Dit is in zekere zin het meest interessante aan de hele abstract–leg dit beter uit: omdat je moet Sigmas toestaat over een arbitrary set (of is het niet beter om een expliciete constraint constructor toe te voegen?), kun je geen generieke generator geven. Dus verwacht je die van de gebruiker.

TODO: Noem heel kort quickchick als alternatief

References

- [1] ALTENKIRCH, T., GHANI, N., HANCOCK, P., MCBRIDE, C., AND MORRIS, P. Indexed containers. *Journal of Functional Programming* 25 (2015).
- [2] CLAESSEN, K., DUREGÅRD, J., AND PÅLKA, M. H. Generating constrained random data with uniform distribution. *Journal of functional programming* 25 (2015).
- [3] CLAESSEN, K., AND HUGHES, J. Quickcheck: a lightweight tool for random testing of haskell programs. *Acm sigplan notices* 46, 4 (2011), 53–64.
- [4] DAGAND, P.-E. *A Cosmology of Datatypes*. PhD thesis, Citeseer, 2013.
- [5] DAGAND, P.-É. The essence of ornaments. *Journal of Functional Programming* 27 (2017).
- [6] PÅLKA, M. H., CLAESSEN, K., RUSSO, A., AND HUGHES, J. Testing an optimising compiler by generating random lambda terms. In *Proceedings of the 6th International Workshop on Automation of Software Test* (2011), ACM, pp. 91–97.
- [7] RUNCIMAN, C., NAYLOR, M., AND LINDBLAD, F. Smallcheck and lazy smallcheck: automatic exhaustive testing for small values. In *Acm sigplan notices* (2008), vol. 44, ACM, pp. 37–48.