

Generating Constrained Test Data using Datatype Generic Programming

Cas van der Rest

Utrecht University

Contents

- Introduction
- · Agda Formalization of 2 type universes
- Regular Types
- · Indexed Descriptions
- · Implementation in Haskell
- · Conclusion

Introduction

Suppose we have the following QuickCheck property

```
propl :: [Int] -> [Int] -> Property
propl xs ys = sorted xs && sorted ys ==> sorted (merge xs ys)
```

What happens when we test this property?

```
sorted xs && sorted ys ==> sorted (merge xs ys)
```

What happens when we test this property?

```
sorted xs && sorted ys ==> sorted (merge xs ys)
```

*** Gave up! Passed only 22 tests; 1000 discarded tests.

The vast majority of generated xs and ys fail the precondition!

What goes wrong here?

What goes wrong here?

 All the generated lists are random, thus unsorted with very high probability

What goes wrong here?

- All the generated lists are random, thus unsorted with very high probability
- $\boldsymbol{\cdot}\;$ A random list that happens to be sorted is likely to be a very short list

We could define a custom generator

```
gen_sorted :: Gen [Int]
gen_sorted = arbitrary >>= return . diff
where diff :: [Int] -> [Int]
    diff [] = []
    diff (x:xs) = x:map (+x) (diff xs)
```

We could define a custom generator

```
gen_sorted :: Gen [Int]
gen_sorted = arbitrary >>= return . diff
where diff :: [Int] -> [Int]
    diff [] = []
    diff (x:xs) = x:map (+x) (diff xs)
```

But what if the we require data with more complex structure (i.e. well-formed programs)

We can represent constrained data as an indexed family

We can represent constrained data as an indexed family

Sorted xs is inhabited if and only if **xs** is a sorted list

We can convert a value of type Sorted xs to a value of type $List \ \mathbb{N}$

```
toList : \forall \{xs\} \rightarrow Sorted \ xs \rightarrow List \ \mathbb{N} toList \{xs\} = xs
```

We can convert a value of type $\textbf{Sorted}\ xs$ to a value of type $\textbf{List}\ \mathbb{N}$

```
toList : \forall {xs} → Sorted xs → List \mathbb{N} toList {xs} \_ = xs
```

If we can generate values of type Sorted xs, we can generate sorted lists!

We can generalize this approach:

We can generalize this approach:

1. Define a suitable indexed type to describe our constrained test data

We can generalize this approach:

- 1. Define a suitable indexed type to describe our constrained test data
- 2. Generate inhabitants of this type

We can generalize this approach:

- 1. Define a suitable indexed type to describe our constrained test data
- 2. Generate inhabitants of this type
- 3. Convert back to the original (non-indexed) datatype

The core contributions of this thesis concern the second step

The core contributions of this thesis concern the second step

More specifically: how can we generically generate values of arbitrary indexed families?

We tackle this question by looking at various type universes, and defining generators for them

Type universes

Each type universe consists of the following elements:

- 1. A datatype **U** describing codes in the universe
- 2. A semantics $[] : U \rightarrow Set$ that maps codes to a type

Type universes

Each type universe consists of the following elements:

- 1. A datatype **U** describing codes in the universe
- 2. A semantics **[_]** : **U** → **Set** that maps codes to a type

Our goal is then to define $deriveGen : (u : U) \rightarrow Gen [u]$

We use an abstract generator type $\textbf{Gen} \;\; \textbf{a} \;\; \textbf{t} \;\; \textbf{i}$ that generates values of indexed types

We use an abstract generator type **Gen a t i** that generates values of indexed types

```
data Gen {i : Set} : Set \rightarrow (i \rightarrow Set) \rightarrow i \rightarrow Set where

Pure : \forall {a t x} \rightarrow a \rightarrow Gen a t x

Or : \forall {a t x} \rightarrow Gen a t x \rightarrow Gen a t x \rightarrow Gen a t x

...

\mu : \forall {t} \rightarrow (x : i) \rightarrow Gen (t x) t x

...

Call : \forall {j t s x} \rightarrow (y : j)

\rightarrow ((y' : j) \rightarrow Gen (s y') s y') \rightarrow Gen (s y) t x
```

We use an abstract generator type **Gen a t i** that generates values of indexed types

```
data Gen {i : Set} : Set → (i → Set) → i → Set where
   Pure : ∀ {a t x} → a → Gen a t x
   Or : ∀ {a t x} → Gen a t x → Gen a t x → Gen a t x
   ...
   μ : ∀ {t} → (x : i) → Gen (t x) t x
   ...
   Call : ∀ {j t s x} → (y : j)
   → ((y' : j) → Gen (s y') s y') → Gen (s y) t x
```

Gen is a deep embedding of the functions exposed by the **Applicative**, **Monad** and **Alternative** typeclasses

We can map **Gen** to any concrete generator type we require

We can map **Gen** to any concrete generator type we require

```
enumerate : \forall {A T} \rightarrow Gen A T \rightarrow Gen T T \rightarrow N \rightarrow List A enumerate g g' zero = [] -- ... enumerate \mu g' (suc n) = enumerate g' g' n
```

We can map Gen to any concrete generator type we require

```
enumerate : \forall {A T} \rightarrow Gen A T \rightarrow Gen T T \rightarrow N \rightarrow List A enumerate g g' zero = [] -- ... enumerate \mu g' (suc n) = enumerate g' g' n
```

We use the enumerative mapping exclusively, but others are possible

Why not skip **Gen** altogether?

```
\begin{array}{lll} \text{nat} & : & (\mathbb{N} \to \text{List } \mathbb{N}) \to \mathbb{N} \to \text{List } \mathbb{N} \\ \\ \text{nat} & \mu = & (|\text{zero}||) \\ \\ & || & (|\text{suc } \mu||) \\ \\ \text{fix} & : & \forall & \{A\} \to (A \to A) \to A \\ \\ \text{fix} & f = f \text{ (fix } f) \\ \end{array}
```

Why not skip **Gen** altogether?

nat :
$$(\mathbb{N} \to \text{List } \mathbb{N}) \to \mathbb{N} \to \text{List } \mathbb{N}$$

nat μ = $($ zero $)$
 $||$ $($ suc μ $)$
fix : \forall $\{A\} \to (A \to A) \to A$
fix f = f $($ fix f $)$

 $\verb"fix" is rejected by the termination checker, using \verb"Gen" circumvents" this issue$

In practice we will allmost never use the constructors of ${\tt Gen}\;\; {\tt a}\;\; {\tt t}\;\; {\tt i}$

In practice we will allmost never use the constructors of $\operatorname{\mathsf{Gen}}\,\,\mathsf{a}\,\,\mathsf{t}\,\,\mathsf{i}$

```
fin : (n : \mathbb{N}) \rightarrow Gen \ (Fin \ n) \ Fin \ n
fin zero = empty
fin (suc \ n) = (zero)
|| (suc \ (\mu \ n))
```

Or desugared:

```
fin : (n : \mathbb{N}) \rightarrow Gen (Fin n) Fin n
fin zero = empty
fin (suc n) = pure zero
 || pure suc <*> (\mu n)
```

Set is isomorphic to the trivial function space $T \rightarrow Set$, so we can generate non-indexed types too

Set is isomorphic to the trivial function space $T \rightarrow Set$, so we can generate non-indexed types too

In practice, this heavily pollutes the code, so we will be a bit liberal with notation

Set is isomorphic to the trivial function space $T \rightarrow Set$, so we can generate non-indexed types too

In practice, this heavily pollutes the code, so we will be a bit liberal with notation

```
nat : Gen \mathbb{N} \mathbb{N}
nat = ( zero )
|| ( suc \mu )
```

instead of

Generator completeness

We want to assert that generators behave correctly

Generator completeness

We want to assert that generators behave correctly

We formulate the following completeness property for this:

```
Complete : \forall {T} \rightarrow Gen T T \rightarrow Set
Complete gen = \forall {x} \rightarrow \exists[ n ] x \in enumerate gen gen n
```

A generator is complete if *all values of the type it produces at some point occur in the enumeration*

Agda Formalization

Universe definition

We start by looking at the universe of Regular Types

Universe definition

We start by looking at the universe of Regular Types

Roughly, this corresponds to algebraic datatypes in Haskell 98

Universe definition

We start by looking at the universe of Regular Types

Roughly, this corresponds to algebraic datatypes in Haskell 98

Examples of regular types include ${\bf Bool}$ and ${\bf List}$

Universe Definition

The universe includes unit types (U), empty types(Z), constant types (K) and recursive positions (I):

```
data Reg : Set where
```

U I Z : Reg

K : Set \rightarrow Reg

Universe Definition

The universe includes unit types (U), empty types(Z), constant types (K) and recursive positions (I):

```
data Reg : Set where
```

U I Z : Reg

K : Set → Reg

Regular types are closed under the product and coproduct operations:

 $_$ $_$ $_$ $_$ $_$ $_$ Reg $_$ Reg $_$ Reg $_$ Reg $_$ Reg

Regular types - Universe Definition

For example: **Bool** is a choice between two nullary constructors:

data Bool : Set where

true : Bool
false : Bool

Regular types - Universe Definition

For example: **Bool** is a choice between two nullary constructors:

data Bool : Set where

true : Bool
false : Bool

Hence, we can describe it as the coproduct of two unit types:

Bool : Reg Bool = $U \oplus U$

Regular types - Semantics

```
The semantics, [\![\ ]\!]: Reg \rightarrow Set \rightarrow Set , map a value of type Reg to a value in Set \rightarrow Set
```

Regular types - Semantics

```
The semantics, [ ]: Reg \rightarrow Set \rightarrow Set, map a value of type Reg to a
value in Set → Set
[ ] : Reg \rightarrow Set \rightarrow Set
[ Z ] r = \bot
■ U
      ] r = T
I I I r = r
\llbracket K X \rrbracket r = X
[ C_1 \otimes C_2 ] r = [ C_1 ] r \times [ C_2 ] r
\llbracket C_1 \oplus C_2 \rrbracket r = \llbracket C_1 \rrbracket r \biguplus \llbracket C_2 \rrbracket r
```

r is the type of recursive positions!

Regular types - Fixpoint operation

We use the following fixpoint operation:

```
data Fix (c : Reg) : Set where
In : [ c ] (Fix c) \rightarrow Fix c
```

Regular types - Fixpoint operation

We use the following fixpoint operation:

```
data Fix (c : Reg) : Set where
  In : [ c ] (Fix c) → Fix c

Fix (U ⊕ U) is isomorphic to Bool.
```

We now aim to define generators from values in Reg

```
deriveGen : (c c' : Reg) \rightarrow \mbox{Gen ([ c ] (Fix c')) ([ c' ] (Fix c'))} \label{eq:condition}
```

We now aim to define generators from values in Reg

```
deriveGen : (c c' : Reg) \rightarrow \mbox{Gen } (\mbox{$ [ \mbox{$ c \mbox{$} ] $ (Fix c')) ( \mbox{$ [ \mbox{$ c' \mbox{$} ] $ (Fix c')) } } \label{eq:constraint}
```

Notice the difference between the type parameters of **Gen**!

```
deriveGen Z c' = empty deriveGen U c' = (|tt|) deriveGen I c' = (|tn|) deriveGen (c_1 \otimes c_2) c' = (|deriveGen c_1 c'|) , (deriveGen c_2 c') (deriveGen c_2 c') (deriveGen c_2 c') (deriveGen c_2 c')
```

What about **K** (constant types)?

```
deriveGen Z c' = empty deriveGen U c' = (|tt||) deriveGen I c' = (|In||\mu|) deriveGen (c_1 \otimes c_2) c' = (|deriveGen||c_1||c'|) , (deriveGen||c_2||c'|) deriveGen (c_1 \oplus c_2) c' = (|inj_1||deriveGen||c_2||c'|) ||(|inj_2||deriveGen||c_2||c'|)
```

26

The semantics of ${\bf K}$ is the type it carries.

We need the programmer's input to generate values of this type

The semantics of **K** is the type it carries.

We need the programmer's input to generate values of this type

How does the programmer supply the required generators?

We define a *metadata structure* that carries additional information about the types stored in a code

We define a *metadata structure* that carries additional information about the types stored in a code

```
data KInfo (P : Set → Set) : Reg → Set where
   Z~ : KInfo P Z
   U~ : KInfo P U
   I~ : KInfo P I
   \otimes \sim : \forall \{c_1 \ c_2\} \rightarrow KInfo P c_1
                                  \rightarrow KInfo P C<sub>2</sub> \rightarrow KInfo P (C<sub>1</sub> \otimes C<sub>2</sub>)
   \oplus \sim : \forall \{c_1 \ c_2\} \rightarrow KInfo P c_1
                                  \rightarrow KInfo P C<sub>2</sub> \rightarrow KInfo P (C<sub>1</sub> \oplus C<sub>2</sub>)
   K\sim : \forall \{S\} \rightarrow PS \rightarrow KInfo P(KS)
```

We parameterise $\tt deriveGen$ over a metadata structure with type $\tt KInfo \ Gen$

```
deriveGen : (c c' : Reg) → KInfo Gen c 
 → Gen ([ c ] (Fix c')) ([ c' ] (Fix c'))
```

We parameterise $\tt deriveGen$ over a metadata structure with type $\tt KInfo \ Gen$

```
deriveGen : (c c' : Reg) → KInfo Gen c 
 → Gen ([ c ] (Fix c')) ([ c' ] (Fix c'))
```

For constant types, ${\tt deriveGen}$ then simply invokes the supplied generator

deriveGen
$$(K_{)}$$
 c' $(K_{}$ g) = Call g

We prove the completeness of **deriveGen** by induction over the input code:

```
complete-thm : \forall {c c' x} \rightarrow 
 \exists[ n ] (x \in enumerate (deriveGen c c') (deriveGen c' c') n)
```

We prove the completeness of **deriveGen** by induction over the input code:

```
complete-thm : \forall {c c' x} \rightarrow \exists[ n ] (x \in enumerate (deriveGen c c') (deriveGen c' c') n)
```

The cases for **U** and **Z** are trivial

```
complete-thm \{U\} = 1 , here complete-thm \{Z\} \{c'\} \{()\}
```

For product and coproduct, we prove that we combine the derived generators in a completeness preserving manner

For product and coproduct, we prove that we combine the derived generators in a completeness preserving manner

This amounts to proving the following lemmas (in pseudocode):

```
Complete g_1 \rightarrow Complete \ g_2
\rightarrow Complete \ (\{ \ inj_1 \ g_1 \ \} \ | \ | \ ( \ inj_2 \ g_2 \ \} )
Complete g_1 \rightarrow Complete \ g_2 \rightarrow Complete \ (\{ \ g_1 \ , \ g_2 \ \} )
```

Recursive positions (I) are slightly more tricky

```
complete-thm {I} {c'} {In x} with complete-thm {c'} {c'} {x} \dots | prf = {!!}
```

Recursive positions (I) are slightly more tricky

```
complete-thm {I} {c'} {In x} with complete-thm {c'} {c'} {x} \dots | prf = {!!}
```

We **must** pattern match on $\mathbf{In}\ \mathbf{x}$, otherwise the recursive call is flagged by the termination checker

We complete this case by proving a lemma of the form:

```
Complete \mu \rightarrow \text{Complete} ( In \mu )
```

For constant types, we parameterize **complete-thm** over a metadata structure containing proofs

 $\mathsf{KInfo}\ (\lambda\ \mathsf{S}\ \rightarrow\ \Sigma[\ \mathsf{g}\ \mathsf{E}\ \mathsf{Gen}\ \mathsf{S}\ \mathsf{S}\]\ \mathsf{Complete}\ \mathsf{g})$

For constant types, we parameterize **complete-thm** over a metadata structure containing proofs

KInfo (
$$\lambda S \rightarrow \Sigma$$
[g \in Gen S S] Complete g)

We then return the proof stored in the metadata structure

Indexed descriptions

Indexed descriptions - Universe definition

The universe of indexed descriptions is largely derived from the universe of regular types

Indexed descriptions - Universe definition

The universe of indexed descriptions is largely derived from the universe of regular types

```
data IDesc (I : Set) : Set where
  `1 : IDesc I
  `var : I → IDesc I
  _`×_ : IDesc I → IDesc I → IDesc I
```

These correspond to U, I and product in the universe of regular types

Indexed descriptions - Universe definition

The regular coproduct is replaced with a generalized version:

```
`\sigma : (n : \mathbb{N}) → (Fin n → IDesc I) → IDesc I
```

Indexed descriptions - Universe definition

The regular coproduct is replaced with a generalized version:

$$^{\circ}\sigma$$
 : (n : N) → (Fin n → IDesc I) → IDesc I

Constant types are replaced with dependent pairs:

$$\Sigma$$
: (S : Set) \rightarrow (S \rightarrow IDesc I) \rightarrow IDesc I

Indexed descriptions - Universe definition

The regular coproduct is replaced with a generalized version:

$$^{\circ}\sigma$$
 : (n : N) → (Fin n → IDesc I) → IDesc I

Constant types are replaced with dependent pairs:

$$\Sigma$$
: (S : Set) \rightarrow (S \rightarrow IDesc I) \rightarrow IDesc I

We denote the empty type with ' σ 0 λ ()

Indexed descriptions - Semantics

The semantic of '1, 'var, and _'x_ are taken (almost) direrectly from the semantics of regular types

Indexed descriptions - Semantics

The semantic of '1, 'var, and _'x_ are taken (almost) direrectly from the semantics of regular types

Both sigma's are interpreted to a dependent pair:

Indexed descriptions - Fixpoint

We can then describe an indexed type using a function of type $I \rightarrow IDesc\ I$.

Indexed descriptions - Fixpoint

We can then describe an indexed type using a function of type $I \rightarrow IDesc\ I$.

The fixpoint operation associated with this universe is:

```
data Fix {I} (\phi : I \rightarrow IDesc I) (i : I) : Set where In : [ \phi i ] (Fix \phi) \rightarrow Fix \phi i
```

Indexed descriptions - Example

Consider a datatype of trees indexed with their size:

```
data STree : \mathbb{N} → Set where

leaf : STree zero

node : \forall {n m} → STree n → STree m

→ STree (suc (n + m))
```

Indexed descriptions - Example

Consider a datatype of trees indexed with their size:

```
data STree : N → Set where
  leaf : STree zero
  node : ∀ {n m} → STree n → STree m
  → STree (suc (n + m))
```

We can use the following indexed description to describe it

```
STree' : \mathbb{N} \to \mathrm{IDesc} \ \mathbb{N}

STree' zero = `1

STree' (suc n) =

`\Sigma \ (\mathbb{N} \times \mathbb{N}) \ \lambda \ \{ \ (m \ , k)

\to \ \Sigma \ (m + k \equiv n) \ \lambda \ \to \ \mathrm{var} \ m \ \mathrm{`} \times \ \mathrm{var} \ k \ \}
```

The generator has the same structure as for regular types

```
deriveGen : \forall {I i} \rightarrow (\delta : IDesc I) \rightarrow (\phi : I \rightarrow IDesc I) \rightarrow Gen ([ \delta ]I (Fix \phi)) (\lambda i \rightarrow [ \phi i ]I (Fix \phi)) i
```

The generator has the same structure as for regular types

```
deriveGen : \forall {I i} \rightarrow (\delta : IDesc I) \rightarrow (\phi : I \rightarrow IDesc I) \rightarrow Gen ([ \delta ]I (Fix \phi)) (\lambda i \rightarrow [ \phi i ]I (Fix \phi)) i
```

The cases for '1, 'var and 'x are also (almost) the same

For the generalized coproduct, we again need to utilize the monadic structure of generators

For the generalized coproduct, we again need to utilize the monadic structure of generators

```
deriveGen (`σ n T) φ = do
  fn ← Call n genFin
  x ← deriveGen (T fn) φ
  pure (fn , x)
```

genFin n generates values of type Fin n

The generalized coproduct is an instantiation of the dependent pair, so we reuse the definition

The generalized coproduct is an instantiation of the dependent pair, so we reuse the definition

```
deriveGen (`\Sigma S T) \phi = do s \leftarrow {!!} x \leftarrow deriveGen (T s) \phi (fm s) pure (s , x)
```

The generalized coproduct is an instantiation of the dependent pair, so we reuse the definition

```
deriveGen (`\Sigma S T) \phi = do s \leftarrow {!!} x \leftarrow deriveGen (T s) \phi (fm s) pure (s , x)
```

How do we get s?

We define a metadata structure:

We define a metadata structure:

```
data IDescM {I : Set} (P : Set → Set) : IDesc I → Set where `var~ : \forall {i} → IDescM P (`var i) `\Sigma~ : \forall {S T} → P S → ((s : S) → IDescM P (T s)) → IDescM P (`\Sigma S T) ...
```

The remaining constructors are handled similar to regular types

We (again) parameterize **deriveGen** over a metadata structure containing generators

```
deriveGen (`\Sigma S T) \phi (`\Sigma~ g mT) = do s \leftarrow Call g x \leftarrow deriveGen (T s) \phi (mT s) pure (s , x)
```

In the case of **STree**, this means that we have to supply a generator that generates pairs of numbers and proofs that their sum is particular number

+-inv : (n :
$$\mathbb{N}$$
) \rightarrow Gen (Σ ($\mathbb{N} \times \mathbb{N}$) λ { (k , m) \rightarrow n \equiv k + m })

In the case of **STree**, this means that we have to supply a generator that generates pairs of numbers and proofs that their sum is particular number

+-inv : (n :
$$\mathbb{N}$$
) \rightarrow Gen (Σ ($\mathbb{N} \times \mathbb{N}$) λ { (k , m) \rightarrow n \equiv k + m })

By using a metadata structure to generate for dependent pairs, we separate the hard parts of generation from the easy parts

In the case of **STree**, this means that we have to supply a generator that generates pairs of numbers and proofs that their sum is particular number

+-inv : (n :
$$\mathbb{N}$$
) \rightarrow Gen (Σ ($\mathbb{N} \times \mathbb{N}$) λ { (k , m) \rightarrow n \equiv k + m })

By using a metadata structure to generate for dependent pairs, we separate the hard parts of generation from the easy parts

A programmer can influence the generation process by supplying different generators

We use the same proof structure as with regular types

```
complete-thm : \forall {\delta \phi x i} \rightarrow 
 \exists[ n ] (x \in enumerate (deriveGen \delta \phi) 
 (\lambda y \rightarrow deriveGen (\phi y) \phi) i n)
```

We use the same proof structure as with regular types

```
complete-thm : \forall {\delta \phi x i} \rightarrow 
 \exists[ n ] (x \in enumerate (deriveGen \delta \phi) 
 (\lambda y \rightarrow deriveGen (\phi y) \phi) i n)
```

enumerate is slightly altered here to accommodate indexed generators

The cases for '1, 'var and 'x follow naturally from the completeness proof for regular types

```
complete-thm {`1} {\phi} {x} = 1 , here complete-thm {`var i} {\phi} {In x} with complete-thm {\phi i} {\phi} {x} ... | prf = {!!} complete-thm {\delta_1 `× \delta_2} {\phi} {x} = {!!}
```

The cases for '1, 'var and 'x follow naturally from the completeness proof for regular types

```
complete-thm {`1} {\phi} {x} = 1 , here complete-thm {`var i} {\phi} {In x} with complete-thm {\phi i} {\phi} {x} ... | prf = {!!} complete-thm {\delta_1 `× \delta_2} {\phi} {x} = {!!}
```

We require (again) additional lemmas of the form:

```
Complete g_1 \to \mathsf{Complete}\ g_2 \to \mathsf{Complete}\ (\ g_1\ ,\ g_2\ ) Complete \mu \to \mathsf{Complete}\ (\ \mathsf{In}\ \mu\ )
```

The generator for dependent pairs is constructed using a monadic bind

The generator for dependent pairs is constructed using a monadic bind Hence, we need to prove an additional lemma about this operation

```
bind-thm :  \forall \ \{g_1 \ g_2 \ A \ B\} \rightarrow \text{Complete} \ g_1 \rightarrow ((x : A) \rightarrow \text{Complete} \ (g_2 \ x))   \rightarrow \text{Complete} \ (g_1 >>= (\lambda \ x \rightarrow g_2 \ x >>= \lambda \ y \rightarrow \text{pure} \ x \ , \ y))
```

To prove completeness for dependent pairs, we can simply invoke this lemma

```
complete-thm {`\Sigma S T} {\phi} = bind-thm {!!} (\lambda x \rightarrow deriveGen (T x) \phi)
```

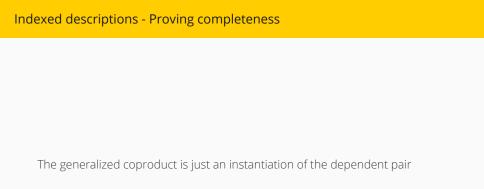
To prove completeness for dependent pairs, we can simply invoke this lemma

```
complete-thm {`\Sigma S T} {\phi} = bind-thm {!!} (\lambda x \rightarrow deriveGen (T x) \phi)
```

The first argument of **bind-thm** is a completeness proof for the user-supplied generator

So we have the user supply this proof

IDescM (λ S \rightarrow Σ [g \in Gen S S] Complete g





The generalized coproduct is just an instantiation of the dependent pair. So we can reuse the proof structure for dependent pairs to prove its completeness

Implementation in Haskell

We make a couple of changes compared to the Agda development:

We make a couple of changes compared to the Agda development:

 The IDesc type gets an extra parameter a, the type that a description describes

We make a couple of changes compared to the Agda development:

- The IDesc type gets an extra parameter a, the type that a description describes
- We represent the generalized coproduct as vector instead of a function

We make a couple of changes compared to the Agda development:

- The IDesc type gets an extra parameter a, the type that a description describes
- We represent the generalized coproduct as vector instead of a function
- We use shallow recursion, meaning that the semantics of recursive position is the associated type a

This means we have no fixpoint combinator!

Universe definition

This all results in the following universe definition

Representing dependent pairs

We choose a more restrictive form of the $\,{}^{\backprime}\Sigma$ combinator, only allowing recursive positions to depend on its first element

Representing dependent pairs

We choose a more restrictive form of the $\,{}^{{}^{{}^{{}}}}\Sigma$ combinator, only allowing recursive positions to depend on its first element

Hence we can change the function $s \rightarrow IDesc \ a \ i$ to a single description $IDesc \ a \ (s \rightarrow i)$.

Representing dependent pairs

We choose a more restrictive form of the $\,{}^{\backprime}\Sigma$ combinator, only allowing recursive positions to depend on its first element

Hence we can change the function $s \to IDesc \ a \ i$ to a single description $IDesc \ a \ (s \to i)$.

Since we use shallow recursion, the semantics of this description is independent of the value of type ${\bf s}$.

Semantics

We describe the semantics in a type family:

We describe the semantics in a type family:

The semantics of the generalized coproduct is just a sum type of all the possible choices

We need a way to express the dependency between the input description, and the type of generated elements

We need a way to express the dependency between the input description, and the type of generated elements

In Agda, we can simply write a Π type: (d : IDesc I) \rightarrow Gen \llbracket c \rrbracket

We need a way to express the dependency between the input description, and the type of generated elements

In Agda, we can simply write a Π type: (d : IDesc I) \rightarrow Gen [c]

In Haskell, we need Singleton types to do this

A singleton type is indexed by some other type, and has exactly one inhabitant for every inhabitant of that type

A singleton type is indexed by some other type, and has exactly one inhabitant for every inhabitant of that type

```
data SNat (n :: Nat) :: * where
   SZero :: SNat Zero
   SSuc :: SNat n -> SNat (Suc n)
```

A singleton type is indexed by some other type, and has exactly one inhabitant for every inhabitant of that type

```
data SNat (n :: Nat) :: * where
   SZero :: SNat Zero
   SSuc :: SNat n -> SNat (Suc n)
inc :: SNat n -> SNat (Suc n)
inc n = SSuc n
```

inc *must* return the successor of its argument, otherwise the typechecker complains!

We define such a singleton type for **IDesc a i** as well:

We define such a singleton type for ${\tt IDesc}$ a ${\tt i}$ as well:

SingIDesc simultaneously acts as a metadata structure, carrying generators for dependent pairs!

With SingIdesc, we can write deriveGen:

```
deriveGen :: SingIDesc d -> Gen (Sem d)
```

For the definition, we follow the Agda implementation.

Using deriveGen

What do we need to use deriveGen?

- 1. A type level description Desc :: i -> IDesc a i
- 2. A singleton instance desc :: Sing i -> SingIDesc (Desc i)
- 3. A conversion function to :: Sing i -> Sem (Desc i) -> a

Consider an expression type:

We'd like to generate well typed expressions (with Type = TNat | TBool)

This comes down to generating values of the following GADT:

```
data Expr (t :: Type) :: * where
AddE :: Expr TNat -> Expr TNat -> Expr TNat
LEQ :: Expr TNat -> Expr TNat -> Expr TBool
ValN :: Nat -> Expr TNat
ValB :: Bool -> Expr TBool
```

We describe this GADT with the following type family ::

```
type family ExprDesc (t :: Type) :: IDesc Expr Type
type instance ExprDesc TNat =
  S2 :+> ( Var TNat :*: Var TNat
          ::: Sigma ('Proxy :: Proxy Nat) One
          ::: VNil )
type instance ExprDesc TBool =
  S2 :+> ( Var TNat :*: Var TNat
          ::: Sigma ('Proxy :: Proxy Bool) One
          ::: VNil )
```

We describe this GADT with the following type family ::

```
type family ExprDesc (t :: Type) :: IDesc Expr Type
type instance ExprDesc TNat =
  S2 :+> ( Var TNat :*: Var TNat
          ::: Sigma ('Proxy :: Proxy Nat) One
          ::: VNil )
type instance ExprDesc TBool =
  S2 :+> ( Var TNat :*: Var TNat
          ::: Sigma ('Proxy :: Proxy Bool) One
          ::: VNil )
And an associated singleton instance: exprDesc :: Sing t ->
SingIDesc (ExprDesc t)
```

Converting back to Expr is then easy:

```
toExpr :: Sing t -> Interpret (ExprDesc t) -> Expr
toExpr STNat (Left (e1 , e2)) = AddE e1 e2
toExpr STNat (Right (n , ()) = ValN n
toExpr STBool (Left (e1 , e2)) = LEQ e1 e2
toExpr STBool (Right (b , ()) = ValB b
```

The definition of toExpr is mostly guided by Haskell's type system

We can now generate well-typed expressions:

```
exprGen :: Sing t -> Gen Expr
exprGen t = toExpr <$> deriveGen (exprDesc t)
```

The elements produced by **exprGen** will all be well-typed expressions.

We can now generate well-typed expressions:

```
exprGen :: Sing t -> Gen Expr
exprGen t = toExpr <$> deriveGen (exprDesc t)
```

The elements produced by exprGen will all be well-typed expressions.

We can use **deriveGen** to generate test data with much richer structure – such as well-typed lambda terms.

Summary

To summarize, we did the following:

- 1. Describe three type universes in Agda, and derive generators from codes in these universes (only two of these discussed here)
- 2. For two of these universes, prove that the generators derived from them are complete
- 3. Implement our development for indexed descriptions in Haskell

We have shown, as a proof of concept, that we can generate arbitrary indexed families

We have shown, as a proof of concept, that we can generate arbitrary indexed families

Of course, this requires that the programmer inputs suitable generators

We have shown, as a proof of concept, that we can generate arbitrary indexed families

Of course, this requires that the programmer inputs suitable generators

With this technique, it is (at least) possible to generate relatively simple well-formed data, such as typed expressions or lambda terms

Future work

Possible avenues for future work include:

- 1. Considering more involved examples, such as polymorphic lambda terms
- 2. Integration with existing testing frameworks
- 3. Applying memoization techniques to the derived generators

Questions?