

# Witnessing (Co)datatypes

Jasmin Christian Blanchette    Andrei Popescu    Dmitriy Traytel

Technische Universität München, Germany  
{blanchet,popescua,traytel}@in.tum.de

## Abstract

Datatypes and codatatypes are useful for specifying and reasoning about (possibly infinite) computational processes. The interactive theorem prover Isabelle/HOL has recently been extended with a definitional package that supports both. Here we describe a complete procedure for deriving nonemptiness witnesses in the general mutually recursive, nested case—nonemptiness being a proviso for introducing new types in higher-order logic. The nonemptiness problem also provides an illuminating case study that shows the package in action, tracing its journey from abstract category theory to hands-on functionality.

**General Terms** Algorithms, Theory, Verification

**Keywords** (Co)datatypes, higher-order logic, interactive theorem proving, Isabelle/HOL, category theory

## 1. Introduction

Interactive theorem provers are becoming increasingly popular as vehicles for formalizing the theory and metatheory of programming languages. Such developments often involve datatypes and codatatypes in various constellations. For example, Lochbihler’s formalization of the Java memory model represents possibly infinite executions using a codatatype [22]. Codatatypes are also useful to capture lazy data structures, such as Haskell’s lists.

Theorem provers based on higher-order logic (HOL), such as HOL4, HOL Light, Isabelle/HOL, and ProofPower–HOL, are traditionally implemented around a trusted inference kernel through which all theorems are generated. Various definitional packages reduce high-level specifications (of types or functions) to primitive inferences; characteristic theorems are derived rather than simply postulated. This reduces the amount of code that must be trusted.

We recently extended Isabelle/HOL with a definitional package for mutually recursive, nested (co)datatypes. While some theorem provers support codatatypes (e.g., Agda, Coq, and PVS), Isabelle is the first to provide a definitional implementation. Our previous paper [32] developed the underlying constructions, adapted from category theory; in this follow-up, we focus on the more practical aspects of the package.

The main such aspect that concerns us here is the generation of nonemptiness witnesses. Types in HOL are required to be nonempty; the “finite stream” specification

$$\text{datatype } \alpha \text{ fstream} = \text{FCons } \alpha (\alpha \text{ fstream})$$

must be rejected because it would lead to an empty datatype. For nonempty types, the package discharges nonemptiness proof obligations by exhibiting existential witnesses.

The nonemptiness problem appears deceptively simple: It is well understood for the mutually recursive datatypes traditionally supported by HOL provers, and for codatatypes the final coalge-

bra is never empty.<sup>1</sup> However, nested recursion via codatatypes and non-free types (e.g., finite sets) complicates the picture, as shown by this attempt at defining infinitely branching trees with finite branches by nested recursion through a codatatype of infinite streams:

$$\begin{aligned} \text{codatatype } \alpha \text{ stream} &= \text{SCons } \alpha (\alpha \text{ stream}) \\ \text{datatype } \alpha \text{ tree} &= \text{Node } \alpha ((\alpha \text{ tree}) \text{ stream}) \end{aligned}$$

The second definition fails: To obtain a witness for  $\alpha \text{ tree}$ , we would need a witness for  $(\alpha \text{ tree}) \text{ stream}$ , and vice versa. Replacing infinite streams with finite lists makes the definition acceptable, because the empty list stops the recursion. Even though final coalgebras are never empty, here the datatype provides a better witness (the empty list) than the codatatype (which requires an  $\alpha \text{ tree}$  to build an  $(\alpha \text{ tree}) \text{ stream}$ ).

Mutual, nested datatype specifications and their nonemptiness witnesses can get complex, as shown by the following commands (aimed at introducing new types with our package):

$$\begin{aligned} \text{datatype } (\alpha, \beta) \text{ tree} &= \\ &\text{Leaf } \beta \mid \text{Branch } ((\alpha + (\alpha, \beta) \text{ tree}) \text{ stream}) \\ \text{codatatype } (\alpha, \beta) \text{ ltree} &= \\ &\text{LNode } \beta ((\alpha + (\alpha, \beta) \text{ ltree}) \text{ stream}) \\ \text{datatype } t_1 &= \\ &T_{11} (((t_1, t_2) \text{ ltree}) \text{ stream}) \mid T_{12} (t_1 \times (t_2 + t_3) \text{ stream}) \\ \text{and } t_2 &= T_2 ((t_1 \times t_2) \text{ list}) \\ \text{and } t_3 &= T_3 ((t_1, (t_3, t_3) \text{ tree}) \text{ tree}) \end{aligned}$$

The definitions are legal but the last group becomes illegal if  $t_2$  is replaced by  $t_3$  in the constructor  $T_{11}$ .

What makes the problem even more interesting is our open-world assumption: the type constructors handled by our package are not syntactically predetermined, in particular, they are not restricted to polynomial functors—the user is allowed to register new type constructors in the package database after checking a few semantic properties. Thus, assume that, in the above mutual definition of the datatypes  $t_i$ , the type constructors  $\text{list}$ ,  $\times$  and  $\text{tree}$  are replaced by others, say,  $\alpha F_1$ ,  $(\alpha, \beta) F_2$ , and  $(\alpha, \beta) F_3$ , respectively, about which we may know nonemptiness information, e.g., that  $\alpha F_1$  is always nonempty, or  $(\alpha, \beta) F_2$  is nonempty whenever  $\alpha$  is. What is the weakest property concerning nonemptiness of the  $F_i$ ’s we need to require in order to have the definition go through? While a lightweight (perhaps syntactic) sufficient condition could in principle prove useful, a condition that is also necessary would allow us to proceed in a *complete* fashion, not rejecting any semantically valid construction.

Our solution, presented in Sect. 4, exploits the package’s abstract, functorial view of types. Each (co)datatype, and more gener-

<sup>1</sup>We ourselves fell prey to this illusion, being led to think that nonemptiness was a minor implementation matter. This explains why the problem is not mentioned in our earlier paper [32], not even as future work.

ally each functor (type constructor) that participates in a definition, carries its own witnesses together with soundness proofs. Operations such as functorial composition, initial algebra, and final coalgebra derive their witnesses from those of their components. Each computational step performed by the package is certified in HOL, and the solution is also shown to be complete.

The other main practical aspect on which we focus is the user interface provided by the package. With the implementation in place (witnesses included), we can complete the journey from abstract category theory to hands-on functionality. This is best illustrated through a case study: a variation of context-free grammars acting on finite sets and their associated possibly infinite derivation trees. The example was carefully chosen to display the support for codatatypes and for nested recursion through non-free types; it also supplies, almost clandestinely, precious building blocks to the non-emptiness proofs of Sect. 4, hence its early placement at Sect. 2.

Finally, Sect. 3 recalls the underlying categorical constructions and connects them to the concrete properties exposed to the user. This section plays a double role, simultaneously acting as an epilogue to Sect. 2 and a prologue to Sect. 4. It also describes the process of deriving customized coinduction principles for codatatypes; this engineering aspect is usually ignored by the category theorist, but it is crucial for the mechanized proof developer. Indeed, the raw categorical structure carried over by the package to perform its type constructions is practically unusable, from both a cultural and a technical point of view: the average Isabelle user should not be expected or required to know category theory; moreover, the outputted theorems should be easy to integrate with the other formal developments, typically not centered around notions such as functoriality or natural transformations, main pillars of our constructions. We show how we have addressed these problems by a compositional approach to customization, taking advantage of the overall compositionality and open-endedness of the collection of types targeted by the package.

The formalization and implementation described in this paper are publicly available [8] and are also part of the latest Isabelle/HOL release [19], where the successful and failing (co)datatype definitions exemplified throughout the paper can now be tested.

## Conventions

We work informally in a mathematical universe  $\mathcal{S}$  of sets but adopt many conventions from higher-order logic and functional programming. Function application is normally written in prefix form without parentheses (e.g.,  $f x y$ ). Sets are ranged over by capital Roman letters ( $A, B, \dots$ ) and Greek letters ( $\alpha, \beta, \dots$ ). For  $n$ -ary functions, we usually prefer the curried form  $f : \alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow \beta$  to the traditional tuple form  $f : \alpha_1 \times \dots \times \alpha_n \rightarrow \beta$  but occasionally pass tuples to curried functions. Polymorphic operators are regarded as families of higher-order constants indexed by sets; thus, the identity function  $\text{id} : \alpha \rightarrow \alpha$  is defined for any set  $\alpha$  and corresponds to a family  $(\text{id}_\alpha)_{\alpha \in \mathcal{S}}$ . Another example is function composition,  $\circ : (\beta \rightarrow \gamma) \rightarrow (\alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \gamma$ .

Operators on sets are normally written in postfix form:  $\alpha \text{ set}$  is the powerset of  $\alpha$ , consisting of sets of elements of  $\alpha$ ;  $\alpha \text{ fset}$  is the set of finite sets over  $\alpha$ . The sum  $\alpha_1 + \alpha_2$  consists of a copy  $\text{Inl } a_1$  of each element  $a_1 : \alpha_1$  and a copy  $\text{Inr } a_2$  of each element  $a_2 : \alpha_2$ . Given  $f_1 : \alpha_1 \rightarrow \beta_1$  and  $f_2 : \alpha_2 \rightarrow \beta_2$ , let  $f_1 \oplus f_2 : \alpha_1 + \alpha_2 \rightarrow \beta_1 + \beta_2$  be the function sending  $\text{Inl } a_1$  to  $\text{Inl } (f_1 a_1)$  and  $\text{Inr } a_2$  to  $\text{Inr } (f_2 a_2)$ . Given  $f : \alpha \rightarrow \beta$ ,  $A \subseteq \alpha$ , and  $B \subseteq \beta$ ,  $\text{image } f A$ , or  $f \cdot A$ , is the image of  $A$  through  $f$ , and  $f^{-1} B$  is the inverse image of  $B$  through  $f$ . The set  $\text{unit}$  contains a single element  $()$ , and  $[n] = \{1, \dots, n\}$ . Prefix and postfix operators bind more tightly than infixes, so that  $\alpha \times \beta \text{ set}$  is read as  $\alpha \times (\beta \text{ set})$  and  $f \cdot g x$  as  $f \cdot (g x)$ .

The notation  $\bar{a}_n$ , or simply  $\bar{a}$ , denotes  $(a_1, \dots, a_n)$ . Given  $\bar{a}_m$  and  $\bar{b}_n$ ,  $(\bar{a}, \bar{b})$  denotes  $(a_1, \dots, a_m, b_1, \dots, b_n)$ . Given  $n$   $m$ -ary func-

tions  $f_1, \dots, f_n$ ,  $\bar{f} \bar{a}$  denotes  $(f_1 \bar{a}, \dots, f_n \bar{a})$ , and similarly  $\bar{a} \bar{F} = (\bar{a} F_1, \dots, \bar{a} F_n)$ . Depending on the context,  $\bar{a}_n F$  either denotes the application of  $F$  to  $\bar{a}$  or merely indicates that  $F$  is an  $n$ -ary set operator.

## 2. The Definitional Package in Action

We introduce the (co)datatype definitional package through a concrete example: derivation trees for a context-free grammar, where we perform the following changes to the usual setting:

- trees are possibly infinite;
- the generated words are not lists, but finite sets.

The formalization of this particular codatatype [8] is the central hub of this paper. First, it is a case study aimed at illustrating some unique characteristics of our package's expressiveness, including the allowance of permutative types such as finite sets inside (co)datatype definitions. Second, it is used to describe the customization process performed by the package. Last but not least, it is employed in the metatheory of arbitrary (co)datatypes needed by the package, namely, in the derivation of nonemptiness witnesses and a proof of soundness and completeness of our algorithm.

We take a few liberties with Isabelle notations to lighten the presentation; in particular, until Sect. 4, we ignore the distinction between types and sets.

### 2.1 Definition of Derivation Trees

We fix a set  $T$  of *terminals* and a set  $N$  of *nonterminals*. The command

```
codatatype dtree = Node (root : N) (cont : (T + dtree) fset)
```

introduces a constructor  $\text{Node} : N \rightarrow (T + \text{dtree}) \text{ fset} \rightarrow \text{dtree}$  and two selectors  $\text{root} : \text{dtree} \rightarrow N$ ,  $\text{cont} : \text{dtree} \rightarrow (T + \text{dtree}) \text{ fset}$ . A tree has the form  $\text{Node } n \text{ as}$ , where  $n$  is a nonterminal (the tree's *root*) and  $\text{as}$  is a finite set of terminals and trees (its *continuation*). The codatatype keyword indicates that this tree formation rule may be applied an infinite number of times. For both datatypes and codatatypes, the package provides basic properties about constructors and selectors, such as injectivity ( $\text{Node } n \text{ as} = \text{Node } n' \text{ as}' \Leftrightarrow n = n' \wedge \text{as} = \text{as}'$ ) and exhaustiveness ( $\forall t. \exists n \text{ as}. t = \text{Node } n \text{ as}$ ).

### Coiterator

For  $\text{dtree}$ , the package also defines the coiterator  $\text{unfold} : (\beta \rightarrow N) \rightarrow (\beta \rightarrow (T + \beta) \text{ fset}) \rightarrow \beta \rightarrow \text{dtree}$  characterized as follows: For all sets  $\beta$ , functions  $r : \beta \rightarrow N$ ,  $c : \beta \rightarrow (T + \beta) \text{ fset}$ , and elements  $b \in \beta$ ,

$$\begin{aligned} \text{root } (\text{unfold } r \ c \ b) &= r \ b \\ \text{cont } (\text{unfold } r \ c \ b) &= (\text{id} \oplus \text{unfold } r \ c) \cdot c \ b \end{aligned}$$

Intuitively, the coiteration contract reads as follows: Given a set  $\beta$ , to define a function  $f : \beta \rightarrow \text{dtree}$  we must indicate how to build a tree for each  $b \in \beta$ . The root is specified by  $r$ , and its continuation is specified corecursively by  $c$ . Formally,  $f = \text{unfold } r \ c$ .

### Structural coinduction

Proofs by coinduction are supported by the structural rule

$$\frac{\forall t_1 t_2. \theta t_1 t_2 \Rightarrow \text{root } t_1 = \text{root } t_2 \wedge \text{fset\_rel } (\text{sum\_rel } (=) \theta) (\text{cont } t_1) (\text{cont } t_2)}{\theta t_1 t_2 \Rightarrow t_1 = t_2}$$

The rule is parameterized by a predicate  $\theta : \text{dtree} \rightarrow \text{dtree} \rightarrow \text{bool}$  that is required to be a bisimulation by the antecedent. The predicate  $\text{fset\_rel } (\text{sum\_rel } (=) \theta)$  is the componentwise extension of  $\theta$  to  $(T + \text{dtree}) \text{ fset}$ . Unfolding the characteristic theorems for  $\text{fset\_rel}$

and `sum_rel` yields the antecedent

$$\begin{aligned} \forall t_1 t_2. \theta t_1 t_2 \Rightarrow & \text{root } t_1 = \text{root } t_2 \wedge \\ & \text{Inl}^-(\text{cont } t_1) = \text{Inl}^-(\text{cont } t_2) \wedge \\ & \forall t'_1 \in \text{Inr}^-(\text{cont } t_1). \exists t'_2 \in \text{Inr}^-(\text{cont } t_2). \theta t'_1 t'_2 \wedge \\ & \forall t'_2 \in \text{Inr}^-(\text{cont } t_2). \exists t'_1 \in \text{Inr}^-(\text{cont } t_1). \theta t'_1 t'_2 \end{aligned}$$

where  $\text{Inl}^-(\text{cont } t)$  is the set of  $t$ 's successor leaves and  $\text{Inr}^-(\text{cont } t)$  is the set of its immediate subtrees. Informally:

If two trees are in relation  $\theta$ , then they have the same root and the same successor leaves and for each immediate subtree of one, there exists an immediate subtree of the other in relation  $\theta$  with it.

The principle effectively states that the equality relation is the largest bisimulation.

## 2.2 Corecursion and Coinduction on Derivation Trees

We start with a simple example of corecursive definition. Let  $+$  :  $\mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$  be a commutative binary operation. We define a parallel composition of trees  $\parallel$  :  $\text{dtree} \times \text{dtree} \rightarrow \text{dtree}$  (written infix), which combines nonterminal-labeled nodes using  $+$  and terminal-labeled leaves using  $\cup$ , corecursively as  $\parallel = \text{unfold } r \ c$ , where  $r$  :  $\text{dtree} \times \text{dtree} \rightarrow \mathbb{N}$  and  $c$  :  $\text{dtree} \times \text{dtree} \rightarrow (\mathbb{T} + \text{dtree} \times \text{dtree})$  fset are

$$\begin{aligned} r(t_1, t_2) &= \text{root } t_1 + \text{root } t_2 \\ c(t_1, t_2) &= \text{Inl} \cdot (\text{Inl}^-(\text{cont } t_1) \cup \text{Inl}^-(\text{cont } t_2)) \cup \\ & \quad \text{Inr} \cdot (\text{Inr}^-(\text{cont } t_1) \times \text{Inr}^-(\text{cont } t_2)) \end{aligned}$$

In accordance with the characteristic coiterator equations,  $\parallel$  is determined by the corecursive equations

$$\begin{aligned} \text{root } (t_1 \parallel t_2) &= r(t_1, t_2) \\ \text{cont } (t_1 \parallel t_2) &= (\text{id} \oplus \parallel) \cdot c(t_1, t_2) \end{aligned}$$

By expanding  $r$  and  $c$  and massaging the second equation, we obtain these characterizations:

$$\begin{aligned} \text{root } (t_1 \parallel t_2) &= \text{root } t_1 + \text{root } t_2 \\ \text{Inl}^-(\text{cont } (t_1 \parallel t_2)) &= \text{Inl}^-(\text{cont } t_1) \cup \text{Inl}^-(\text{cont } t_2) \\ \text{Inr}^-(\text{cont } (t_1 \parallel t_2)) &= \parallel \cdot (\text{Inr}^-(\text{cont } t_1) \times \text{Inr}^-(\text{cont } t_2)) \\ &= \{t'_1 \parallel t'_2 \mid t'_1 \in \text{Inr}^-(\text{cont } t_1) \wedge \\ & \quad t'_2 \in \text{Inr}^-(\text{cont } t_2)\} \end{aligned}$$

The equations specify the root, the successor leaves, and the immediate subtrees of the tree  $t_1 \parallel t_2$ , in that order. To prove that  $\parallel$  is commutative, we first formulate the goal in the format expected by the structural coinduction rule,  $\theta t t' \Rightarrow t = t'$ , with  $\theta t t' \Leftrightarrow \exists t_1 t_2. t = t_1 \parallel t_2 \wedge t' = t_2 \parallel t_1$ . Structural coinduction yields the following subgoals, which are easy to discharge using the equations for  $\parallel$  and commutativity of  $+$  and  $\cup$ :

$$\begin{aligned} \text{root } (t_1 \parallel t_2) &= \text{root } (t_2 \parallel t_1) \\ \text{Inl}^-(\text{cont } (t_1 \parallel t_2)) &= \text{Inl}^-(\text{cont } (t_2 \parallel t_1)) \\ \forall t \in \text{Inr}^-(\text{cont } (t_1 \parallel t_2)). \exists t' \in \text{Inr}^-(\text{cont } (t_2 \parallel t_1)). \theta t t' \\ \forall t' \in \text{Inr}^-(\text{cont } (t_2 \parallel t_1)). \exists t \in \text{Inr}^-(\text{cont } (t_1 \parallel t_2)). \theta t t' \end{aligned}$$

## 2.3 A Variation of Context-Free Grammars

We continue with a larger example that illustrates the coinductive infrastructure offered by the package. The example is a variation of the notion of context-free grammar, acting on finite sets instead of sequences. We assume that the previously fixed sets  $\mathbb{T}$  and  $\mathbb{N}$ , of terminals and nonterminals, are finite and that we are given a set of *productions*  $\mathbb{P} \subseteq \mathbb{N} \times (\mathbb{T} + \mathbb{N})$  fset. The triple  $\text{Gr} = (\mathbb{T}, \mathbb{N}, \mathbb{P})$  forms a (*set*) *grammar*, which is fixed for the rest of this section. Both finite and infinite derivation trees are of interest. The codatatype

`dtree` provides the right universe for defining well-formed trees as a coinductive predicate.

Fixpoint (or Knaster–Tarski) (co)induction is provided in Isabelle/HOL by a separate package [26]. Fixpoint induction relies on the minimality of a predicate (the least fixpoint); dually, fixpoint coinduction relies on maximality (the greatest fixpoint). It is well-known that datatypes interact well with definitions by fixpoint induction. For codatypes, both fixpoint induction and coinduction play an important role—the former to express safety properties, the latter to express liveness.

We illustrate fixpoint (co)induction on `dtree`. *Well-formed* derivation trees for  $\text{Gr}$  are defined coinductively as the greatest predicate  $\text{wf} : \text{dtree} \rightarrow \text{bool}$  such that, for all  $t \in \text{dtree}$ ,

$$\begin{aligned} \text{wf } t \Leftrightarrow & (\text{root } t, (\text{id} \oplus \text{root}) \cdot \text{cont } t) \in \mathbb{P} \wedge \\ & \text{root is injective on } \text{Inr}^-(\text{cont } t) \wedge \\ & \forall t' \in \text{Inr}^-(\text{cont } t). \text{wf } t' \end{aligned}$$

Intuitively, each nonterminal node of a well-formed derivation tree  $t$  represents a production. This is achieved formally by three conditions: (1) the root of  $t$  forms a production together with the terminals constituting its successor leaves and the roots of its immediate subtrees; (2) no two immediate subtrees of  $t$  have the same root; (3) properties 1 and 2 also hold for the immediate subtrees of  $t$ . The definition's coinductive nature ensures that these properties hold for arbitrarily deep subtrees of  $t$ , even if  $t$  has infinite depth.

In contrast to wellformedness, the notions of subtree, frontier (the set of terminals appearing in a tree), and interior (the set of nonterminals appearing in a tree) require inductive definitions. The *subtree* relation is defined as the least predicate  $\text{subtr} : \text{dtree} \rightarrow \text{dtree} \rightarrow \text{bool}$  such that  $\text{subtr } t t' \Leftrightarrow t = t' \vee (\exists t''. \text{subtr } t t'' \wedge \text{Inr } t'' \in \text{cont } t')$  holds for all  $t, t' \in \text{dtree}$ .<sup>2</sup> A coinductive definition would be unsuitable here:  $\text{subtr } t t'$  would be true for any infinite tree  $t'$  by virtue of the second disjunct (the step). We write  $\text{Subtr } t$  for the set of subtrees of  $t$ . The frontier  $\text{Fr} : \text{dtree} \rightarrow \mathbb{T}$  set and interior  $\text{ltr} : \text{dtree} \rightarrow \mathbb{N}$  set of a tree are defined similarly. The language generated by the grammar  $\text{Gr}$  from a nonterminal  $n \in \mathbb{N}$  (using possibly infinite derivation trees) is defined as  $\mathcal{L}_{\text{Gr}}(n) = \{\text{Fr } t \mid \text{wf } t \wedge \text{root } t = n\}$ .

The choice of the permutative structure of finite sets instead of lists to represent words has important consequences for the computation of the generated language, where regular trees suffice. We elaborate on this next, taking the opportunity to illustrate corecursion.

A derivation tree is *regular* if each subtree is uniquely determined by its root. Formally, we define *regular*  $t$  as the existence of a function  $f : \mathbb{N} \rightarrow \text{Subtr } t$  such that  $\forall t' \in \text{Subtr } t. f(\text{root } t') = t'$ . The regularly generated language of a nonterminal is defined as  $\mathcal{L}_{\text{Gr}}^r(n) = \{\text{Fr } t \mid \text{wf } t \wedge \text{root } t = n \wedge \text{regular } t\}$ .

Given a possibly nonregular derivation tree  $t_0$ , a *regular cut* of  $t_0$  is a regular tree  $\text{rcut } t_0$  such that  $\text{Fr}(\text{rcut } t_0) \subseteq \text{Fr } t_0$ . Here is one way to perform the cut:

1. Choose a subtree of  $t_0$  for each interior node  $n \in \text{ltr } t_0$  via a function  $\text{pick} : \text{ltr } t_0 \rightarrow \text{Subtr } t_0$  with  $\forall n \in \text{ltr } t_0. \text{root}(\text{pick } n) = n$ .
2. Traverse  $t_0$  and replace each subtree with  $\text{root } n$  with  $\text{pick } n$ . The replacement should be performed hereditarily, i.e., also in the emerging subtree  $\text{pick } n$ .

This replacement task is elegantly achieved by the corecursive function  $\text{H} : \text{ltr } t_0 \rightarrow \text{dtree}$  defined as  $\text{unfold } r \ c$ , where  $r : \text{ltr } t_0 \rightarrow \mathbb{T}$  and  $c : \text{ltr } t_0 \rightarrow \mathbb{T} + (\text{ltr } t_0)$  fset are specified as follows:  $r n = n$  and  $c n = (\text{id} \oplus \text{root}) \cdot \text{cont}(\text{pick } n)$ .  $\text{H}$  is therefore characterized by the

<sup>2</sup> Inductive predicates are often specified by their introduction rules—here,  $\text{subtr } t t$  and  $\text{subtr } t t'' \wedge \text{Inr } t'' \in \text{cont } t' \Rightarrow \text{subtr } t t'$ . The formulation as an equivalence is convenient for its uniform treatment of the coinductive case.

corecursive equations

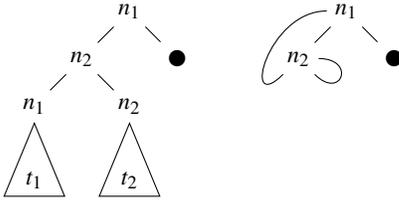
$$\begin{aligned} \text{root } (H n) &= n \\ \text{cont } (H n) &= (\text{id} \oplus (H \circ \text{root})) \cdot \text{cont } (\text{pick } n) \end{aligned}$$

It is not hard to prove the following by fixpoint coinduction:

**Theorem 1.** For all  $n \in \text{ltr } t_0$ ,  $H n$  is regular and  $\text{Fr } (H n) \subseteq \text{Fr } t_0$ . Moreover,  $H n$  is well-formed provided  $t_0$  is well-formed.

(The proof is given in Appendix D.) Therefore we take  $\text{rcut } t_0$  to be  $H(\text{root } t_0)$ .

Figure 1 illustrates a derivation tree and its simplest regular cut. The bullet denotes a terminal, and  $t_1$  and  $t_2$  are arbitrary trees with roots  $n_1$  and  $n_2$ . The loops in the right-hand tree denote infinite trees that are their own subtrees.



**Fig. 1.** A derivation tree (left) and its simplest regular cut (right)

### 3. The Journey from Category Theory

The previous section's definitions rest on characteristic theorems provided by Isabelle's (co)datatype package—they are derived from underlying constructions adapted from category theory. The fundamental concept is that of bounded natural functors, a well-behaved class of functors with additional structure. The theory is an essential preliminary to the nonemptiness witness computation of Sect. 4 and puts the case study of Sect. 2 on more solid ground.

#### 3.1 Functors and Functor Operations

We consider various operators  $F$  on sets, such as sums, products, etc., which we call *set constructors*. We are interested in set constructors that are *functors* on the category of sets and functions, meaning that they are equipped with an action on morphisms commuting with identities and composition. This action is a polymorphic constant  $\text{Fmap} : (\alpha_1 \rightarrow \beta_1) \rightarrow \dots \rightarrow (\alpha_n \rightarrow \beta_n) \rightarrow \bar{\alpha} F \rightarrow \bar{\beta} F$  satisfying  $\text{Fmap } \text{id}^n = \text{id}$  and  $\text{Fmap } (g_1 \circ f_1) \dots (g_n \circ f_n) = \text{Fmap } \bar{g} \circ \text{Fmap } \bar{f}$ . Formally, functors are pairs  $(F, \text{Fmap})$ . We let  $\text{Func}_n$  be the collection of  $n$ -ary functions. Let us review basic examples.

##### Identity functor (ID, id)

The identity functor maps any set and any function to itself.

##### $(n, \alpha)$ -constant functor ( $C_{n,\alpha}$ , $C_{\text{map}_{n,\alpha}}$ )

The  $(n, \alpha)$ -constant functor  $(C_{n,\alpha}, C_{\text{map}_{n,\alpha}})$  is the  $n$ -ary functor consisting of the set constructor  $\bar{\beta} C_{n,\alpha} = \alpha$  and the action  $C_{\text{map}_{n,\alpha}} f_1 \dots f_n = \text{id}$ . We write  $C_\alpha$  for  $C_{1,\alpha}$ .

##### Sum functor (+, $\oplus$ )

The involved operators have already been described in Sect. 1.

##### Product functor ( $\times$ , $\otimes$ )

Let  $\text{fst} : \alpha_1 \times \alpha_2 \rightarrow \alpha_1$  and  $\text{snd} : \alpha_1 \times \alpha_2 \rightarrow \alpha_2$  denote the two projection functions. Given  $f_1 : \alpha \rightarrow \beta_1$  and  $f_2 : \alpha \rightarrow \beta_2$ , let  $\langle f_1, f_2 \rangle : \alpha \rightarrow \beta_1 \times \beta_2$  be the function  $\lambda a. (f_1 a, f_2 a)$ . Given  $f_1 : \alpha_1 \rightarrow \beta_1$  and  $f_2 : \alpha_2 \rightarrow \beta_2$ , let  $f_1 \otimes f_2 : \alpha_1 \times \alpha_2 \rightarrow \beta_1 \times \beta_2$  be  $\langle f_1 \circ \text{fst}, f_2 \circ \text{snd} \rangle$ .

##### $\alpha$ -Function space functor ( $\text{func}_\alpha$ , $\text{comp}_\alpha$ )

Given a set  $\alpha$ , let  $\beta \text{func}_\alpha = \alpha \rightarrow \beta$ . For all  $g : \beta \rightarrow \gamma$ , let  $\text{comp}_\alpha g : \beta \text{func}_\alpha \rightarrow \gamma \text{func}_\alpha$  be  $\text{comp}_\alpha g f = g \circ f$ .

##### Powerset functor (set, image)

For all  $f : \alpha \rightarrow \beta$ , the function  $\text{image } f : \alpha \text{ set} \rightarrow \beta \text{ set}$  sends each subset  $A$  of  $\alpha$  to the image of  $A$  through the function  $f : \alpha \rightarrow \beta$ .

##### Bounded $k$ -powerset functor ( $\text{set}_k$ , image)

Given a cardinal  $k$ , for all sets  $\alpha$ , the set  $\alpha \text{set}_k$  carves out from  $\alpha \text{ set}$  only those sets of cardinality less than  $k$ . The finite powerset functor  $\text{fset}$  corresponds to  $\text{set}_{\aleph_0}$ .

Functors can be composed to form more complex functors. Composition requires the functors  $F_j$  to take the same type arguments  $\bar{\alpha}$  in the same order. The auxiliary operations of permutation and lifting, together with the identity and  $(n, \alpha)$ -constant functors (Sect. 3.1), make it possible to compose functors freely.

##### Composition

Given  $\bar{\alpha} F_j$  for  $j \in [n]$  and  $\bar{\beta}_n G$ , the *functor composition*  $G \circ \bar{F}$  is defined on objects as  $(\bar{\alpha} \bar{F}) G$  and similarly on morphisms.

##### Permutation

Given  $F \in \text{Func}_n$  and  $i, j \in [n]$  with  $i < j$ , the  $(i, j)$ -permutation of  $F$ , written  $F^{(i,j)} \in \text{Func}_n$ , is defined on objects as  $\bar{\alpha} F^{(i,j)} = (\alpha_1, \dots, \alpha_{i-1}, \alpha_j, \alpha_{i+1}, \dots, \alpha_{j-1}, \alpha_i, \alpha_{j+1}, \dots, \alpha_n) F$  and similarly on morphisms.

##### Lifting

Given  $F \in \text{Func}_n$ , the *lifting* of  $F$ , written  $F \uparrow \in \text{Func}_{n+1}$ , is defined on objects as  $(\bar{\alpha}_n, \alpha_{n+1}) F \uparrow = \bar{\alpha}_n F$  and similarly on morphisms. In other words,  $F \uparrow$  is obtained from  $F$  by adding a superfluous argument.

Datatypes are defined by taking the initial algebra of a set of functors and codatatype by taking the final coalgebra. Both operations are partial.

##### Initial algebra

Given  $n$   $(m+n)$ -ary functors  $(\bar{\alpha}_m, \bar{\beta}_n) F_j$ , their *(mutual) initial algebra* consists of  $n$   $m$ -ary functors  $\bar{\alpha} \text{IF}_j$  that satisfy the isomorphism  $\bar{\alpha} \text{IF}_j \cong (\bar{\alpha}, \bar{\alpha} \text{IF}) F_j$  minimally. (The variables  $\bar{\alpha}$  are the passive parameters, and  $\bar{\beta}$  are the fixpoint variables.) The functors  $\text{IF}_j$  are characterized by

- $n$  polymorphic *folding bijections* (constructors)  $\text{ctor}_j : (\bar{\alpha}, \bar{\alpha} \text{IF}) F_j \rightarrow \bar{\alpha} \text{IF}_j$  and
- $n$  polymorphic *iterators*  $\text{fold}_j : (\prod_{k \in [n]} (\bar{\alpha}, \bar{\beta}) F_k \rightarrow \beta_k) \rightarrow \bar{\alpha} \text{IF}_j \rightarrow \beta_j$

subject to the following properties (for all  $j \in [n]$ ):

- Iteration equations:  $\text{fold}_j \bar{s} \circ \text{ctor}_j = s_j \circ \text{Fmap } \text{id}^m (\text{fold } \bar{s})$ .
- Unique characterization of iterators: Given  $\bar{\beta}$  and  $\bar{s}$ , the only functions  $f_j : \bar{\alpha} \text{IF}_j \rightarrow \beta_j$  satisfying  $f_j \circ \text{ctor}_j = s_j \circ \text{Fmap } \text{id}^m \bar{f}$  are  $\text{fold}_j \bar{s}$ .

The functorial actions  $\text{IFmap}_j$  for  $\text{IF}_j$  are defined by iteration in a standard way.

##### Final coalgebra

The final coalgebra operation is categorically dual to initial algebra. Given  $n$   $(m+n)$ -ary functors  $(\bar{\alpha}_m, \bar{\beta}_n) F_j$ , their *(mutual) final*

coalgebra consists of  $n$   $m$ -ary functors  $\bar{\alpha} JF_j$  that satisfy the isomorphism  $\bar{\alpha} JF_j \cong (\bar{\alpha}, \bar{\alpha} JF) F_j$  maximally. The functors  $JF_j$  are characterized by

- $n$  polymorphic *unfolding bijections* (destructors)  $\text{dtr}_j : \bar{\alpha} JF_j \rightarrow (\bar{\alpha}, \bar{\alpha} JF) F_j$  and
- $n$  polymorphic *coiterators*  $\text{unfold}_j : (\prod_{k \in [n]} \beta_k \rightarrow (\bar{\alpha}, \bar{\beta}) F_k) \rightarrow \beta_j \rightarrow \bar{\alpha} JF_j$

subject to the following properties:

- Coiteration equations:  $\text{dtr}_j \circ \text{unfold}_j \bar{s} = \text{Fmap id}^m (\overline{\text{unfold}} \bar{s}) \circ s_j$ .
- Unique characterization of coiterators: Given  $\bar{\beta}$  and  $\bar{s}$ , the only functions  $f_j : \beta_j \rightarrow \bar{\alpha} JF_j$  satisfying  $\text{dtr}_j \circ f_j = \text{Fmap id}^m \bar{f} \circ s_j$  are  $\text{unfold}_j \bar{s}$ .

The functorial actions  $JF\text{map}_j$  for  $JF_j$  are defined by coiteration in the standard way.

The (co)datatype package is based on a class  $\mathcal{B}$  of functors, called *bounded natural functors (BNFs)* [32], such that  $\mathcal{B}$  contains all the basic functors listed in Sect. 3.1 except for unbounded power-set, and  $\mathcal{B}$  supports, and is closed under, the above operations.

For the derivation tree example of Sect. 2.1, the input BNF to the initial algebra operation is  $\text{pre\_dtree} = (\times) \circ (\text{C}_N, \text{fset}) \circ ((+) \circ (\text{C}_T, \text{ID}))$ . In the sequel, we prefer the more readable notation  $\alpha \text{pre\_dtree} = N \times (T + \alpha) \text{fset}$ .

### 3.2 Open-Endedness

Unlike the (co)datatype specification mechanisms of most theorem provers (such as the previous one of Isabelle and the current ones of the other HOL-based provers, as well as those of PVS and Coq), here the involved types are not syntactically predetermined by a fixed grammar.  $\mathcal{B}$  does include the class of polynomial functors, but is open-ended in that users may register further functors as members of  $\mathcal{B}$ .

The registration process takes place as follows. The user provides a type constructor  $F$  and its associated BNF infrastructure (in the form of polymorphic HOL constants), including the  $\text{Fmap}$  functorial action on objects. Then the user verifies the BNF properties, e.g., that  $(F, \text{Fmap})$  is indeed a functor. (The full list of BNF-defining properties can be found in [32].) After this, the new BNF is integrated and can appear nested in future (co)datatype definitions. For instance, Isabelle users have recently introduced the BNFs  $\alpha \text{bag}$ , of bags (i.e., bags with finite-multiplicity elements) over  $\alpha$ , and  $\alpha \text{list}_5$ , of lists of at most five elements. (More details are given in Appendix E.)

### 3.3 Abstract Structural (Co)induction

Besides closure under functor operations, another important question for theorem proving is how to state induction and coinduction abstractly, irrespective of the shape of the functor. We know how to state induction on lists, or trees, but how about on initial algebras of arbitrary functors? The answer is based on enriching the structure of functors  $\bar{\alpha}_n F$  with additional data: For each  $i \in [n]$ , BNFs must provide a natural transformation  $\text{Fset}^i : \bar{\alpha} F \rightarrow \alpha_i \text{set}$  that gives, for  $x \in \bar{\alpha} F$ , the set of  $\alpha_i$ -atoms that take part in  $x$ . For example, if  $(\alpha_1, \alpha_2) F = \alpha_1 \times \alpha_2$ , then  $\text{Fset}^1(a_1, a_2) = \{a_1\}$  and  $\text{Fset}^2(a_1, a_2) = \{a_2\}$ ; if  $\alpha F = \alpha \text{list}$  (the list functor, obtained as minimal solution to  $\beta \cong \text{unit} + \alpha \times \beta$ ), then  $\text{Fset}$  applied to a list  $x$

gives all the elements appearing in  $x$ .<sup>3</sup> As usual, if  $F$  is unary, we omit the superscript from  $\text{Fset}$ <sup>1</sup>.

Given  $j \in [n]$ , the elements of  $\text{Fset}_j^{m+k} x$  (for  $k \in [n]$ ) are the recursive components of  $\text{ctor}_j x$ . (Note that here subscripts select functors  $F_j$  in the tuple  $\bar{F}$ ; then, as previously agreed, superscripts select  $\text{Fset}$  operators for different arguments of  $F_j$ .) Using this insight, the induction principle can be expressed abstractly for the mutual initial algebra  $\bar{F}$  of functors  $\bar{F}$  as follows for sets  $\bar{\alpha}$  and predicates  $\varphi_j : \bar{\alpha} JF_j \rightarrow \text{bool}$ :

$$\frac{\bigwedge_{j=1}^n \forall x \in (\bar{\alpha}, \bar{\alpha} JF) F_j. (\bigwedge_{k=1}^n \forall b \in \text{Fset}_j^{m+k} x. \varphi_k b) \Rightarrow \varphi_j (\text{ctor}_j x)}{\bigwedge_{j=1}^n \forall b \in \bar{\alpha} JF_j. \varphi_j b}$$

For lists, this gives the abstract principle

$$\frac{\forall x \in \text{unit} + \alpha \times \alpha \text{list}. (\forall b \in \text{Fset}^2 x. \varphi b) \Rightarrow \varphi (\text{ctor } x)}{\forall b \in \alpha \text{list}. \varphi b}$$

which, by taking  $\text{Nil} = \text{ctor} (\text{Inl } ())$  and  $\text{Cons } a b = \text{ctor} (\text{Inr } (a, b))$ , can be recast into the familiar rule

$$\frac{\varphi \text{Nil} \quad \forall a \in \alpha. \forall b \in \alpha \text{list}. \varphi b \Rightarrow \varphi (\text{Cons } a b)}{\forall b \in \alpha \text{list}. \varphi b}$$

Moving to coinduction, we need a further well-known assumption [29]: that our functors preserve weak pullbacks, which allows us to organize them as *relators* [28]. For a functor  $\bar{\alpha}_n F$ , we lift its action  $\text{Fmap} : (\alpha_1 \rightarrow \beta_1) \rightarrow \dots \rightarrow (\alpha_n \rightarrow \beta_n) \rightarrow \bar{\alpha} F \rightarrow \bar{\beta} F$  on functions to an action  $\text{Frel} : (\alpha_1 \rightarrow \beta_1 \rightarrow \text{bool}) \rightarrow \dots \rightarrow (\alpha_n \rightarrow \beta_n \rightarrow \text{bool}) \rightarrow (\bar{\alpha} F \rightarrow \bar{\beta} F \rightarrow \text{bool})$ , the *relator*, defined as follows:

$$\text{Frel } \bar{\varphi} x y \Leftrightarrow \exists z. \text{Fmap } \bar{\text{fst}} z = x \wedge \text{Fmap } \bar{\text{snd}} z = y \wedge \bigwedge_{i=1}^n \forall (a, b) \in \text{Fset}^i z. \varphi_i a b$$

Structural coinduction can be stated abstractly as

$$\frac{\bigwedge_{j=1}^n \forall a b \in (\bar{\alpha}, \bar{\alpha} JF) F_j. \theta_j a b \Rightarrow \text{Frel}_j (=) \bar{\theta} (\text{dtr}_j a) (\text{dtr}_j b)}{\bigwedge_{j=1}^n \forall a b. \theta_j a b \Rightarrow a = b}$$

for sets  $\bar{\alpha}_n$  and binary predicates  $\theta_j \in \bar{\alpha} JF_j \rightarrow \bar{\alpha} JF_j \rightarrow \text{bool}$  [28]. The antecedent ensures that  $\theta$  is a bisimulation.

### 3.4 Completing the Journey: From Abstract to Concrete

Given the definition of  $\text{dtree}$  (Sect. 2.1), the package first composes the intermediate BNF  $\alpha \text{pre\_dtree} = N \times (T + \alpha) \text{fset}$  from the constants  $N$  and  $T$ , identity, sum, product, and finite set. Then it constructs the final coalgebra  $\text{dtree} (= JF)$  from  $\text{pre\_dtree} (= F)$ .

The unfolding bijection  $\text{dtr} : \text{dtree} \rightarrow \text{dtree } \text{pre\_dtree}$  is decomposed in two selectors:  $\text{root} = \text{fst} \circ \text{dtr}$  and  $\text{cont} = \text{snd} \circ \text{dtr}$ . The constructor  $\text{Node}$  is defined as the inverse of the unfolding bijection. The basic properties of constructors and selectors are derived from those of sums and products.

#### Customizing corecursion and coinduction

The abstract coiteration principle described in Sect. 3.1 relies on a coiterator  $\text{unfold} : (\beta \rightarrow \beta \text{pre\_dtree}) \rightarrow \beta \rightarrow \text{dtree}$  such that  $\text{dtr} \circ \text{unfold } s = \text{pre\_dtree\_map} (\text{unfold } s) \circ s$ . Writing  $s$  as  $\langle r, c \rangle$  for  $r : \beta \rightarrow N$  and  $c : \beta \rightarrow (T + \alpha) \text{fset}$  and recasting the equation in pointful form, we obtain  $\text{dtr} (\text{unfold } \langle r, c \rangle b) = \text{pre\_dtree\_map} (\text{unfold } s) (r b, c b)$ . This can be further improved by unfolding the definition of  $\text{pre\_dtree\_map}$ , expressing  $\text{dtr}$  as  $\langle \text{root}, \text{cont} \rangle$ , and splitting the result into a pair of equations:  $\text{root} (\text{unfold } \langle r, c \rangle b) = r b$  and

<sup>3</sup> Our  $\text{Fset}$  has similarities with Pierce's notion of support from his account of (co)inductive types [27] and with Abel and Altenkirch's urelement relation from their framework for strong normalization [2]. The main distinguishing feature of our notion is that it additionally takes categorical structure into consideration [32].

$\text{cont } (\text{unfold } \langle r, c \rangle b) = (\text{id} \oplus \text{unfold } \langle r, c \rangle) \cdot c b$ . The coiteration rule of Sect. 3.1 emerges by replacing  $\text{unfold}$  with the curried  $\text{unfold}' : (\beta \rightarrow \mathbb{N}) \rightarrow (\beta \rightarrow (\top + \beta)) \text{fset} \rightarrow \beta \rightarrow \text{dtree}$  defined as  $\text{unfold}' r c = \text{unfold } \langle r, c \rangle$ .

Customizing the abstract structural coinduction principle from Sect. 3.3 requires more work. From the abstract, functor-agnostic notion of bisimulation, our goal is to derive concrete coinduction rules that avoid any references to the intermediate functor  $F$  (e.g.,  $\text{pre\_dtree}$ ). Category-theoretical approaches typically appeal to ad hoc proofs that the abstract formulation of bisimulation, in terms of the graph of  $F\text{map}$ , is equivalent to the traditional bisimulation principle for various interesting concrete cases (e.g., [29, Ex. 2.1]). However, in a theorem prover, we prefer a uniform, automatic customization that yields the expected concrete instances.

The key is to employ an alternative compositional characterization. To this end, each standard basic BNF provides an ad hoc characterization. For example:

- Product:  $\text{prod\_rel } \theta_1 \theta_2 (a_1, a_2) (b_1, b_2) \Leftrightarrow \theta a_1 b_1 \wedge \theta a_2 b_2$
- Bounded  $k$ -Powerset:  $\text{set\_rel}_k \theta A B \Leftrightarrow (\forall a \in A. \exists b \in B. \theta a b) \wedge (\forall b \in B. \exists a \in A. \theta a b)$

The same can be achieved for any potential user-defined BNF: The user is prompted to provide an ad hoc characterization of the relator. (If the user provides none, then the general-purpose definition in terms of  $F\text{map}$  is kept.) E.g., for bags, a convenient such characterization proceeds inductively by the following clauses, where  $\uplus$  denotes bag union (by summing the elements' multiplicities):

$$\text{bag\_rel } \theta \emptyset \emptyset \qquad \frac{\text{bag\_rel } \theta A B \quad \theta a b}{\text{bag\_rel } \theta (A \uplus \{a\}) (B \uplus \{b\})}$$

Moreover, each BNF operation derives a characterization of the resulting relator in terms of those of the components.

- Composition:  $\text{Hrel } \bar{\theta} = \text{Grel } (\overline{\text{Frel}} \bar{\theta})$
- Initial algebra:  $\text{IFrel}_j \bar{\theta} (\text{ctor}_j a) (\text{ctor}_j b) \Leftrightarrow \text{Frel}_j \bar{\theta} (\overline{\text{IFrel}} \bar{\theta}) a b$
- Final coalgebra:  $\text{JFrel}_j \bar{\theta} a b \Leftrightarrow \text{Frel}_j \bar{\theta} (\overline{\text{JFrel}} \bar{\theta}) (\text{dctor}_j a) (\text{dctor}_j b)$

The last two equivalences hold in a fixpoint-(co)inductive way: For the initial algebra, the  $\text{IFrel}_j$ 's are the least predicates (mutually) satisfying the equivalence; for the final coalgebra, the  $\text{JFrel}_j$ 's are the greatest predicates satisfying the equivalence.

For  $\alpha$   $\text{pre\_dtree} = \mathbb{N} \times (\top + \alpha)$   $\text{fset}$  and its final coalgebra  $\text{dtree}$ , the relators are characterized by the following equivalences:

$$\begin{aligned} \text{pre\_dtree\_rel } \theta a b &\Leftrightarrow \\ \text{prod\_rel } (\text{C}_{\mathbb{N}\text{-rel}} \theta) (\text{fset\_rel } (\text{sum\_rel } (\text{C}_{\top\text{-rel}} \theta) (\text{ID\_rel } \theta))) a b & \\ \text{dtree\_rel } a b &\Leftrightarrow \text{pre\_dtree\_rel } \text{dtree\_rel } (\text{dctor } a) (\text{dctor } b) \end{aligned}$$

The final formulation, as presented in Sect. 2.1, emerges by substituting  $\langle \text{root}, \text{cont} \rangle$  for  $\text{dctor}$  and unfolding the characteristic equations for the involved basic BNFs.

## 4. Computing Nonemptiness Witnesses

In the previous two sections, we referred to the codatatype  $\text{dtree}$  and other collections of elements as *sets*, postponing the discussion of an important aspect of our package. While for most purposes sets and types can be identified in an abstract treatment of HOL, types have the additional restriction that they may not be empty. The main primitive way to define custom types in HOL is to specify from an existing type  $\alpha$  a nonempty subset  $A : \alpha$   $\text{set}$  that is isomorphic to the desired type. Hence, to register a collection of elements as a HOL type it is necessary to prove it nonempty.

Mutual datatype definitions are a particular case of the above situation, with the additional requirement that the nonemptiness proof

should be performed automatically by the package. In the context of our package, we need to produce the relevant nonemptiness proofs taking into consideration arbitrary combinations of datatypes, codatatypes and user-defined BNFs.

A first temptation to tackle the problem is to follow the traditional approach of HOL datatype packages [6, 10]: Try to unfold all the definitions of the involved nested datatypes, inlining them as additional components of the mutual definition, until only sums of products remain, and then perform a reachability analysis. At a closer inspection, this approach turns out problematic in our framework for several reasons:

- Due to open-endedness, there is no fixed set of basic types.
- Delving into nested types requires re-proving nonemptiness facts, which is extremely inefficient.
- It is not clear how to unfold datatypes nested in codatatypes or vice versa.

Note that, counting on everything being eventually reducible to the fixed situation of sums of products, the traditional approach worries about nonemptiness only at datatype-definition time. Here, we look for a prophylactic solution instead, trying to prepare the BNFs in advance for future nonemptiness checks involving them. To this end, we ask the following: Given a mutual datatype definition involving several  $n$ -ary BNFs, what is the relevant information we need to know about their nonemptiness *without knowing how they look like* (hence, with no option to “delve” into them)? To answer this, we use a generalization of pointed types [17, 21], maintaining witnesses that assert conditional nonemptiness for combinations of arguments. We first present the solution by examples.

### 4.1 Examples

We start with the very simple cases of products and sums. For  $\alpha \times \beta$ , the proof is as follows: Assuming  $\alpha \neq \emptyset$  and  $\beta \neq \emptyset$ , we construct the witness  $(a, b) \in \alpha \times \beta$  for some  $a \in \alpha$  and  $b \in \beta$ . For  $\alpha + \beta$ , two proofs are possible: Assuming  $\alpha \neq \emptyset$ , we can construct  $\text{Inl } a$  for some  $a \in \alpha$ ; alternatively, assuming  $\beta \neq \emptyset$ , we can construct  $\text{Inr } b$  for some  $b \in \beta$ .

To each BNF  $\bar{\alpha} F$ , we associate a set of witnesses, each of the form  $\text{Fwit} : \alpha_{i_1} \rightarrow \dots \rightarrow \alpha_{i_k} \rightarrow \bar{\alpha} F$  for a subset  $\{i_1, \dots, i_k\} \subseteq [n]$ . From a witness, we can construct a set-theoretic proof by following the signature of the witness (in the spirit of the Curry–Howard correspondence). Accordingly,  $\text{Inr} : \beta \rightarrow \alpha + \beta$  can be read as the following contract: Given a proof that  $\beta$  is nonempty,  $\text{Inr}$  yields a proof that  $\alpha + \beta$  is nonempty.

When BNFs are composed, so are their witnesses. Thus, the two possible witnesses for the list-defining functor  $(\alpha, \beta)$   $\text{pre\_list} = \text{unit} + \alpha \times \beta$  are  $\text{pre\_list\_wit}_1 = \text{Inl } ()$  and  $\text{pre\_list\_wit}_2 a b = \text{Inr } (a, b)$ . The first witness subsumes the second one, because it unconditionally shows the collection nonempty, regardless of the potential emptiness of  $\alpha$  and  $\beta$ . From this witness, we obtain the unconditional witness  $\text{list\_ctor } \text{pre\_list\_wit}_1$  (i.e.,  $\text{Nil}$ ) for  $\alpha$   $\text{list}$ .

Because they can store infinite objects, codatatypes are never empty. Compare the following:

$$\begin{aligned} \text{datatype } \alpha \text{ fstream} &= \text{FSCons } \alpha (\alpha \text{ fstream}) \\ \text{codatatype } \alpha \text{ stream} &= \text{SCons } \alpha (\alpha \text{ stream}) \end{aligned}$$

The datatype definition fails because the best witness has a circular signature:  $\alpha \rightarrow \alpha \text{ fstream} \rightarrow \alpha \text{ fstream}$ . In contrast, the codatatype definition succeeds and produces the witness  $(\lambda a. \mu s. \text{SCons } a s) : \alpha \rightarrow \alpha \text{ stream}$ , namely the (unique) stream  $s$  such that  $s = \text{SCons } a s$  for a given  $a \in \alpha$ . This stream is easy to define by coiteration.

Let us look at a pair of examples involving nesting:

$$\begin{aligned} \text{datatype } (\alpha, \beta) \text{ tree} &= \\ \text{Leaf } \beta \mid \text{Branch } ((\alpha + (\alpha, \beta) \text{ tree}) \text{ stream}) & \end{aligned}$$

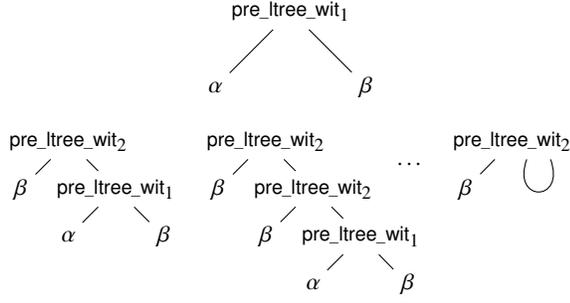


Fig. 2. Derivation trees for ltree witnesses

codatatype  $(\alpha, \beta)$  ltree =  
 LNode  $\beta$  (( $\alpha + (\alpha, \beta)$  ltree) stream)

In the tree definition, the two constructors hide a sum BNF, giving us some flexibility. For the Leaf constructor, all we need is a witness  $b \in \beta$ , from which we construct Leaf  $b$ . For Branch, we can choose the left-hand side of the nested  $+$ , completely dodging the recursive right-hand side: From a witness  $a \in \alpha$ , we construct Branch ( $\mu s$ . SCons (Inl  $a$ )  $s$ ).

For the ltree functor, the two arguments to LNode are hiding a product, so the ltree-defining functor is  $(\alpha, \beta, \gamma)$  pre\_ltree =  $\beta \times (\alpha + \gamma)$  stream with  $\gamma$  representing the corecursive component. Composition yields two witnesses for pre\_ltree:

- $\text{pre\_ltree\_wit}_1 \ a \ b = (b, \mu s. \text{SCons} (\text{Inl } a) \ s)$
- $\text{pre\_ltree\_wit}_2 \ b \ c = (b, \mu s. \text{SCons} (\text{Inr } c) \ s)$

These can serve to build infinitely many witnesses for ltree. Figure 2 enumerates the possible combinations, starting with  $\text{pre\_ltree\_wit}_1$ . This witness requires only the non-corecursive components  $\alpha$  and  $\beta$  being nonempty, and hence immediately yields a witness  $\text{pre\_ltree\_wit}_1 : \alpha \rightarrow \beta \rightarrow (\alpha, \beta)$  ltree (by applying the constructor LNode). The second witness  $\text{pre\_ltree\_wit}_2$  requires both  $\beta$  and the corecursive component  $\gamma$  to be nonempty; it effectively “consumes” another ltree witness through  $\gamma$ . The consumed witness can again be either  $\text{pre\_ltree\_wit}_1$  or  $\text{pre\_ltree\_wit}_2$ , and so on. At the limit,  $\text{pre\_ltree\_wit}_2$  is used infinitely often. The corresponding witness  $\text{pre\_ltree\_wit}_2 : \beta \rightarrow (\alpha, \beta)$  ltree can be defined by coiteration as  $\lambda b. \mu t. \text{pre\_ltree\_wit}_2 \ b \ t$ . It subsumes  $\text{pre\_ltree\_wit}_1$  and all the other finite witnesses. Were ltree to be defined as a datatype instead of a codatatype,  $\text{pre\_ltree\_wit}_1$  would be its best witness.

## 4.2 A General Solution

The nonemptiness problem for an  $n$ -ary set constructor  $F$  and a set of indices  $I \subseteq [n]$  can be stated as follows: Does it hold that, for all sets  $\bar{\alpha}_n, \bar{\alpha} \ F \neq \emptyset$  whenever  $\forall i \in I. \alpha_i \neq \emptyset$ ? We call  $F$   $I$ -witnessed if the above question has a positive answer. E.g., set sum ( $+$ ) is  $\{1\}$ -,  $\{2\}$ -, and  $\{1, 2\}$ -witnessed; set product ( $\times$ ) is  $\{1, 2\}$ -witnessed; and  $\alpha$  list is  $\emptyset$ -witnessed.

We are led to the following notion of soundness. Given an  $n$ -ary functor  $F$ , a set  $\mathcal{I} \subseteq [n]$  set is (*witness-sound*) for  $F$  if  $F$  is  $I$ -witnessed for all  $I \in \mathcal{I}$ .

Now, when is such a set  $\mathcal{I}$  also complete, in that it covers all witnesses? To answer this, first note that, if  $I_1 \subseteq I_2$ , then  $I_1$ -witnesshood implies  $I_2$ -witnesshood. Therefore, we are interested in retaining the witnesses completely only up to inclusion of sets of indices. We thus call a set  $\mathcal{I} \subseteq [n]$  set

- (*witness-complete*) for  $F$  if for all  $J \subseteq [n]$  such that  $F$  is  $J$ -witnessed, there exists  $I \in \mathcal{I}$  such that  $I \subseteq J$ ;
- (*witness-perfect*) for  $F$  if it is both sound and complete.

Here are perfect sets  $\mathcal{I}_F$  for some basic BNFs:

- Identity:  $\mathcal{I}_{\alpha \text{ ID}} = \{\{\alpha\}\}$
- Constant:  $\mathcal{I}_{C_{n, \alpha}} = \{\emptyset\}$  ( $\alpha \neq \emptyset$ )
- Sum:  $\mathcal{I}_{\alpha + \beta} = \{\{\alpha\}, \{\beta\}\}$
- Product:  $\mathcal{I}_{\alpha \times \beta} = \{\{\alpha, \beta\}\}$
- Function space:  $\mathcal{I}_{\beta \text{ func}_\alpha} = \{\{\beta\}\}$  ( $\alpha \neq \emptyset$ )
- Bounded  $k$ -Powerset:  $\mathcal{I}_{\alpha \text{ set}_k} = \{\emptyset\}$
- Bag:  $\mathcal{I}_{\alpha \text{ bag}} = \{\emptyset\}$

Parameters  $\alpha_j$  are identified with their indices  $j$  to improve readability.

We need to maintain perfect sets across BNF operations. Let us start with composition, permutation, and lifting.

**Theorem 2.** Let  $H = G \circ \bar{F}_n$ , where  $G \in \text{Func}_n$  has a perfect set  $\mathcal{I}$  and each  $F_j \in \text{Func}_m$  has a perfect set  $\mathcal{I}_j$ . Then  $\{\bigcup_{j \in J} I \mid J \in \mathcal{I} \wedge (I_j)_{j \in J} \in \prod_{j \in J} \mathcal{I}_j\}$  is a perfect set for  $H$ .

**Theorem 3.** Let  $\mathcal{I} \subseteq [n]$  be a perfect set for  $F$ . Then  $\mathcal{I}$  is also perfect for  $F^{(i, j)}$  and  $F \uparrow$ .

Theorems 2 and 3 hold not only for functors but also for plain set constructors (with a further cardinality-monotonicity assumption needed for the completeness part of Theorem 2). The most interesting cases are the genuinely functorial ones of initial algebras and final coalgebras, which we discuss next.

Witnesses for the initial algebra and the final coalgebra will be essentially obtained by repeated compositions of the witnesses of the involved BNFs and the folding bijections, inductively in one case and coinductively in the other. The derivation trees from Sect. 2.3 turn out to be perfectly suited for recording the combinatorics of these compositions, in such a way that not only soundness, but also completeness yield easily.

For the rest of this subsection, we fix  $n$  ( $m+n$ )-ary functors  $\bar{F} F_j$  and assume each  $F_j$  has a perfect set  $\mathcal{K}_j$ .

We construct a (set) grammar  $\text{Gr} = (\text{T}, \text{N}, \text{P})$  as in Sect. 2.3 with  $\text{T} = [m]$ ,  $\text{N} = [n]$ , and  $\text{P} = \{(j, \text{cp}(K)) \mid K \in \mathcal{K}_j\}$ , where, for each  $K \subseteq [m+n]$ ,  $\text{cp}(K)$  is its copy to  $[m] + [n]$  defined as  $\text{Inl} \cdot ([m] \cap K) \cup \text{Inr} \cdot \{k \in [n] \mid m+k \in K\}$ .

Here is the idea behind this construction. A mutual datatype definition as above introduces  $n$  isomorphisms:

$$\begin{aligned} \bar{\alpha} \ \text{IF}_1 &\cong (\bar{\alpha}, \bar{\alpha} \ \text{IF}_1, \dots, \bar{\alpha} \ \text{IF}_n) \ F_1 \\ &\vdots \\ \bar{\alpha} \ \text{IF}_n &\cong (\bar{\alpha}, \bar{\alpha} \ \text{IF}_1, \dots, \bar{\alpha} \ \text{IF}_n) \ F_n \end{aligned}$$

We are searching for conditions guaranteeing nonemptiness of the  $\text{IF}_j$ 's. To this end, we walk these isomorphisms from left to right, reducing nonemptiness of  $\bar{\alpha} \ \text{IF}_j$  to that of  $(\bar{\alpha}, \bar{\alpha} \ \text{IF}_1, \dots, \bar{\alpha} \ \text{IF}_n) \ F_j$ . Moreover, nonemptiness of the latter can be reduced to nonemptiness of some  $\alpha_{i_1}, \dots, \alpha_{i_p}$  and some  $\bar{\alpha} \ \text{IF}_{j_1}, \dots, \bar{\alpha} \ \text{IF}_{j_q}$ , along a witness for  $F_j$  of the form  $\{i_1, \dots, i_p\} \cup \{m+j_1, \dots, m+j_q\}$ . This yields a grammar production  $j \rightarrow \{\text{Inl } i_1, \dots, \text{Inl } i_p\} \cup \{\text{Inr } j_1, \dots, \text{Inr } j_q\}$ , where the  $i_k$ 's are terminals and the  $j_l$ 's are, like  $j$ , nonterminals. The ultimate goal is to eventually obtain reductions of the nonemptiness of  $\bar{\alpha} \ \text{IF}_j$  to that of components of  $\bar{\alpha}$  alone, i.e., to terminals—this precisely corresponds to derivations in the grammar of terminal sets. It should be intuitively clear that by considering finite derivations we obtain sound witnesses for  $\text{IF}_j$ . We shall actually prove more:

- for initial algebras, finite derivations are also witness-complete;
- for final coalgebras (replacing  $\bar{\text{IF}}$  with  $\bar{\text{JF}}$ ), accepting possibly infinite derivations is still sound, and also becomes complete.

**Theorem 4.** Assume that the (mutual) final coalgebra of  $\bar{F}$  exists and consists of  $n$   $m$ -ary functors  $\bar{\alpha}_m \text{JF}_j$  (as in Sect. 3.1). Then  $\mathcal{L}_{\text{Gr}}^r(j)$  is a perfect set for  $\text{JF}_j$  for  $j \in [n]$ .

To prove soundness, we define a nonemptiness witness to  $\bar{\alpha} \text{JF}_j$  corecursively (by abstract  $\bar{\text{JF}}$ -corecursion). More interestingly: To prove completeness, we define a function to  $\text{dtree}$  corecursively (by concrete tree corecursion), obtaining a derivation tree, from which we then cut a regular derivation tree via Theorem 1.

*Proof sketch:* Let  $j_0 \in [n]$ . We first show that  $\mathcal{L}_{\text{Gr}}^r(j_0)$  is sound. Let  $t_0$  be a well-formed regular derivation tree with root  $j_0$ . We need to prove that  $F_{j_0}$  is  $\text{Fr } t_0$ -witnessed. For this, we fix  $\bar{\alpha}_m$  such that  $\forall i \in \text{Fr } t_0, \alpha_i \neq \emptyset$ , and aim to show that  $\bar{\alpha} \text{JF}_{j_0} \neq \emptyset$ .

For each  $j \in \text{ltr } t_0$ , let  $t_j$  be the corresponding subtree of  $t_0$ . (It is well-defined, since  $t_0$  is regular.) Note that  $t_0 = t_{j_0}$ . For each  $K$  such that  $(j, \text{cp}(K)) \in \text{P}$ , since  $K \in \mathcal{K}_j$  and  $\mathcal{K}_j$  is sound for  $F_j$ , we obtain a  $K$ -witness for  $F_j$ , i.e., a function  $w_{j,K} : (\gamma_k)_{k \in K} \rightarrow \bar{\gamma} F_j$  (polymorphic in  $\bar{\gamma}$ ).

Let  $\beta_n$  be defined as  $\beta_j = \text{unit}$  if  $j \in \text{ltr } t_0$  and  $= \emptyset$  otherwise. We build a coalgebra structure on  $\beta$ ,  $(s_j : \beta_j \rightarrow (\bar{\alpha}, \beta) F_j)_{j \in [n]}$ , as follows: If  $j \notin \text{ltr } t_0$ ,  $s_j$  is the unique function from  $\emptyset$ . If  $j \in \text{ltr } t_0$ , then  $s_j() = w_{j,K}(a_i)_{i \in K \cap [m]}()^{K \cap [m+1, m+n]}$ , where  $\text{cp}(K)$  is the right-hand side of the top production of  $t_j$ , i.e.,  $(\text{id} \oplus \text{root}) \cdot \text{cont } t_j$ . Now, for each  $j \in \text{ltr } t_0$ ,  $\text{unfold}_j \bar{s} : \text{unit} \rightarrow \bar{\alpha} \text{JF}_j$  ensures the nonemptiness of  $\bar{\alpha} \text{JF}_j$ . In particular,  $\bar{\alpha} \text{JF}_{j_0} \neq \emptyset$ .

We now show that  $\mathcal{L}_{\text{Gr}}^r(j_0)$  is complete. Let  $I \subseteq [m]$  such that  $\text{JF}_{j_0}$  is  $I$ -witnessed. We need to find  $I_1 \in \mathcal{L}_{\text{Gr}}^r(j_0)$  such that  $I_1 \subseteq I$ . Let  $\bar{\alpha}_m$  be defined as  $\alpha_i = \text{unit}$  if  $i \in I$  and  $= \emptyset$  otherwise. Let  $J = \{j, \bar{\alpha} F_j \neq \emptyset\}$ . We define  $c : J \rightarrow ([m] + J)$  fset by  $c j = \text{cp}(K_j)$ , where  $K_j$  is such that  $(j, \text{cp}(K_j)) \in \text{P}$  and  $K_j \subseteq I \cup \{m + j, j \in J\}$ .

Let now  $g : J \rightarrow \text{dtree}$  be  $\text{unfold id } c$ . Thus, for all  $j \in J$ ,  $\text{root}(g j) = j$  and  $\text{cont}(g j) = (\text{id} \oplus g) \cdot c j = \text{Inl} \cdot (K_j \cap I) \cup \text{Inr} \cdot \{g j, m + j \in K_j\}$ . Taking  $t_0 = g j_0$  and using Theorem 1, we obtain the regular well-formed tree  $t_1$  such that  $\text{Fr } t_1 \subseteq \text{Fr } t_0 \subseteq I$ . Hence  $\text{Fr } t_1$  is the desired  $I_1$ .  $\square$

The above completeness proof provides an example of self-application of codatatypes: A specific codatatype, of infinite derivation trees, figures in the metatheory of general codatatypes. And this may well be unavoidable: While for soundness the regular trees are replaceable by some equivalent (finite) inductive items, it is not clear how completeness could be proved without first appealing to arbitrary infinite derivation trees and then cutting them down to regular trees.

An analogous result holds for initial algebras. For each  $i \in \mathbb{N}$ , let  $\mathcal{L}_{\text{Gr}}^r(i)$  be the language generated by  $i$  by means of regular finite Gr derivation trees. Since  $\mathbb{N}$  is finite, these can be described more directly as trees for which every nonterminal path has no repetitions.

**Theorem 5.** Assume that the (mutual) initial algebra of  $\bar{F}$  exists and consists of  $n$   $m$ -ary functors  $\bar{\alpha}_m \text{IF}_j$  (as in Sect. 3.1). Then  $\mathcal{L}_{\text{Gr}}^r(j)$  is a perfect set for  $\text{IF}_j$  for  $j \in [n]$ .

Let us see how Theorems 2–5 can be combined in establishing or refuting nonemptiness for some of our motivating examples from Sect. 4.1 and the introduction. These (and other) examples have been checked with our package [8]:

- $\mathcal{S}_{(\alpha, \beta) \text{pre\_list}} = \{\emptyset\}$  by Th. 2;  $\mathcal{S}_{\alpha \text{list}} = \{\emptyset\}$  by Th. 5
- $\mathcal{S}_{(\alpha, \beta) \text{pre\_fstream}} = \{\{\alpha, \beta\}\}$ ;  
 $\mathcal{S}_{\alpha \text{fstream}} = \emptyset$  by Th. 5 (i.e.  $\alpha \text{fstream}$  is empty)
- $\mathcal{S}_{(\alpha, \beta) \text{pre\_stream}} = \{\{\alpha, \beta\}\}$ ;  $\mathcal{S}_{\alpha \text{stream}} = \{\{\alpha\}\}$  by Th. 4
- $\mathcal{S}_{(\alpha, \beta) \text{pre\_x}} = \{\{\beta\}\}$  and  $\mathcal{S}_{(\alpha, \beta) \text{pre\_y}} = \{\emptyset\}$ ;  
 $\mathcal{S}_x = \{\emptyset\}$  and  $\mathcal{S}_y = \{\emptyset\}$  by Th. 5

- $\mathcal{S}_{(\alpha, \beta, \gamma) \text{pre\_ltree}} = \{\{\alpha, \beta\}, \{\beta, \gamma\}\}$  by Th. 2;  
 $\mathcal{S}_{(\alpha, \beta) \text{ltree}} = \{\{\beta\}\}$  by Th. 4
- $\mathcal{S}_{(\alpha, \beta, \gamma) \text{pre\_t}_1} = \{\{\beta\}, \{\alpha, \gamma\}\}$ ,  $\mathcal{S}_{(\alpha, \beta, \gamma) \text{pre\_t}_2} = \{\emptyset\}$ , and  
 $\mathcal{S}_{(\alpha, \beta, \gamma) \text{pre\_t}_3} = \{\{\alpha\}, \{\gamma\}\}$  by Th. 2;  
 $\mathcal{S}_{t_i} = \{\emptyset\}$  by Th. 5

Since we have maintained perfect sets throughout all the BNF operations, we obtain the following central result:

**Theorem 6.** Any BNF built from BNFs endowed with perfect sets of witnesses (in particular all basic BNFs discussed in this paper) by repeated applications of the composition, initial algebra, and final coalgebra operations has a perfect set defined as indicated in Theorems 2–5.

Consequently, a procedure implementing Theorems 2–5 will preserve enough nonemptiness witnesses to ensure that all datatype specifications describing nonempty types are accepted. The next subsection presents such a procedure.

### 4.3 Computational Aspects

Theorem 4 reduces the computation of perfect sets for final coalgebras to that of  $\mathcal{L}_{\text{Gr}}^r(n)$ . Our use of infinite regular trees in the definition of  $\mathcal{L}_{\text{Gr}}^r(n)$  had the advantage of allowing a simple proof of soundness, and the only natural proof of completeness we could think of, relating the coinductive nature of arbitrary mutual codatatypes with that of infinite trees. However, from the computational point of view, regular infinite trees are of course excessive—we next show that  $\mathcal{L}_{\text{Gr}}^r(n)$  can be computed directly via recursive equations not involving derivation trees.

First, we relativize the notion of frontier to that of “frontier through  $ns$ ,”  $\text{Fr } ns$   $t$ , consisting of the leaves of  $t$  that are accessible through paths of terminals belonging to  $ns \subseteq \mathbb{N}$ . We also define the corresponding  $ns$ -restricted regularly generated language  $\mathcal{L}_{\text{Gr}}^r ns$ .

Recall that our derivation trees from Sect. 2 produce not usual words (defined as lists), but finite sets—in what follows, by “word” we mean “finite set of terminals”. We can think of a generated word as being more precise than another provided the former is a subword (subset) of the latter. This leads us to defining, for languages (sets of words), the notions of word-inclusion subsumption<sup>4</sup>,  $\leq$ , by

$$L \leq L' \text{ iff } \forall w \in L. \exists w' \in L'. w' \subseteq w$$

and equivalence,  $\equiv$ , by

$$L \equiv L' \text{ iff } L \leq L' \text{ and } L' \leq L$$

It is easy to see that any set  $\equiv$ -equivalent to a perfect set is again perfect. Note also that Theorem 1 implies  $\mathcal{L}_{\text{Gr}}^r(n) \equiv \mathcal{L}_{\text{Gr}}(n)$ , which qualifies regular trees as a generated-language optimization of arbitrary trees.

We compute  $\mathcal{L}_{\text{Gr}}^r ns$  up to word-inclusion equivalence  $\equiv$  by recursively applying the available productions whose source nonterminals are in  $ns$ , removing each time the expanded nonterminal from  $ns$ . Thus, provided  $n$  is in the allowed set  $ns$ ,  $\mathcal{L}_{\text{Gr}}^r ns$  calls  $\mathcal{L}_{\text{Gr}}^r ns' n'$  recursively with  $ns' = ns \setminus \{n'\}$  for each nonterminal  $n'$  in the chosen production from  $n$ , and so on, until the current node is no longer in the decumulator set  $ns$ . Formally:

**Theorem 7.** For all  $ns \subseteq \mathbb{N}$  and  $n \in \mathbb{N}$ ,  $\mathcal{L}_{\text{Gr}}^r ns$   $n \equiv$

$$\begin{cases} \{\emptyset\} & \text{if } n \notin ns \\ \{\text{Inl}^- ss \cup \bigcup_{n' \in \text{Inr}^- ss} K_{n'} \mid (n, ss) \in \text{P} \wedge K \in \prod_{n' \in \text{Inr}^- ss} \mathcal{L}_{\text{Gr}}^r (ns \setminus \{n\}) n'\} & \text{otherwise} \end{cases}$$

Theorem 7 provides an alternative, recursive definition of  $\mathcal{L}_{\text{Gr}}^r ns$   $n$ . The definition terminates because the argument  $ns$  is fi-

<sup>4</sup>This is in effect the Smyth preorder extension [31] of the subword relation.

nite and decreases strictly in the recursive case—in fact, this shows that the height of the recursive call stack is bounded by the number of nonterminals, which for our intended application translates to the number of simultaneously introduced codatatypes.

Here is how this recursion operates on the tree example. We have  $T = \{\alpha, \beta\}$ ,  $N = \{\gamma\}$ , and  $P = \{p_1, p_2\}$ , where  $p_1 = (\gamma, \{\text{Inl } \alpha, \text{Inl } \beta\})$  and  $p_2 = (\gamma, \{\text{Inl } \beta, \text{Inr } \gamma\})$ . Note that

- $\text{Inl}^- ss = \{\alpha, \beta\}$  and  $\text{Inr}^- ss = \emptyset$  for  $(n, ss) = p_1$
- $\text{Inl}^- ss = \{\beta\}$  and  $\text{Inr}^- ss = \{\gamma\}$  for  $(n, ss) = p_2$

The computation has one single recursive call, yielding

$$\begin{aligned} \mathcal{L}_{\text{Gr}}^r \gamma &= \mathcal{L}_{\text{Gr}}^r \{\gamma\} \gamma \\ &\equiv \{\{\alpha, \beta\} \cup \emptyset\} \cup \\ &\quad \{\{\beta\} \cup \bigcup_{n' \in \{\gamma\}} K_{n'} \mid K \in \prod_{n' \in \{\gamma\}} \mathcal{L}_{\text{Gr}}^r \emptyset n'\} \\ &= \{\{\alpha, \beta\}\} \cup \{\{\beta\} \cup K_\gamma \mid K_\gamma \in \mathcal{L}_{\text{Gr}}^r \emptyset \gamma\} \\ &= \{\{\alpha, \beta\}\} \cup \{\{\beta\} \cup \emptyset\} \\ &= \{\{\alpha, \beta\}, \{\beta\}\} \\ &\equiv \{\{\beta\}\} \end{aligned}$$

For datatypes, the computation of  $\mathcal{L}_{\text{Gr}}^{\text{rf}}$  is achieved analogously to Theorem 7, by defining  $\mathcal{L}_{\text{Gr}}^{\text{rf}} ns n$  as a generalization of  $\mathcal{L}_{\text{Gr}}^{\text{rf}} n$ .

**Theorem 8.** The statement of Theorem 7 still holds if we substitute  $\mathcal{L}_{\text{Gr}}^{\text{rf}}$  for  $\mathcal{L}_{\text{Gr}}^r$  and  $\emptyset$  for  $\{\emptyset\}$ .

## 5. Implementation in Isabelle/ML

The package maintains nonemptiness information to be prepared for producing nonemptiness proofs upon datatype definitions. The equations from Theorems 7 and 8 involve only executable operations over finite sets of numbers, sums, and products. Since the descriptions of Theorems 2 and 3 are also executable, the implementation task emerges clearly: Store a perfect set along each basic BNF and have each BNF operation compute witnesses from those of its operands.

However, as it stands,  $I$ -witnesshood cannot be expressed in HOL because types are always nonempty: How can we state that  $(\alpha, \beta)$  tree  $\neq \emptyset$  conditionally on  $\alpha \neq \emptyset$  or  $\beta \neq \emptyset$ , in the context of  $\alpha$  and  $\beta$  being assumed nonempty in the first place? The solution is to work not with operators  $\bar{\alpha} F$  on HOL types directly but rather with their *internalization* to sets, expressed as a polymorphic function

$$\text{Fin} : \alpha_1 \text{ set} \rightarrow \dots \rightarrow \alpha_n \text{ set} \rightarrow (\bar{\alpha} F) \text{ set}$$

defined as  $\text{Fin } \bar{A} = \{x. \forall i \in [n]. \text{Fset}^i x \subseteq A_i\}$ .  $I$ -witnesshood becomes  $(\forall i \in I. A_i \neq \emptyset) \Rightarrow \text{Fin } \bar{A} \neq \emptyset$ .

For each  $n$ -ary BNF  $F$ , the package stores a set of sets  $\mathcal{S}$  of numbers in  $[n]$  (the perfect set) and, for each set  $I \in \mathcal{S}$ , a polymorphic constant  $w_I : (\alpha_i)_{i \in I} \rightarrow \bar{\alpha} F$  and an equivalent formulation of  $I$ -witnesshood:  $\forall i \in I. \text{Fset}^i (w_I (a_j)_{j \in I}) \neq \emptyset$ .

Due to the logic's restricted expressiveness, we cannot prove the theorems presented in this paper in their most general form for arbitrary functors and have the package instantiate them for specific functors. Instead, the package proves the theorems dynamically for the specific functors involved in the datatype definitions. Only the soundness part of the theorems is needed. To paraphrase Krauss and Nipkow [20], completeness belongs to the realm of metatheory and is not required to obtain actual nonemptiness proofs—it does let you sleep better though, by ensuring that the employed criterion is as precise as it can be.

A HOL definitional package has to bear the burden of both *computing* HOL terms and *certifying* the computation, i.e., ensuring that certain terms are theorems. The combinatorial computation of witnessing sets of indices described in Theorems 7 and 8 would be expensive if performed through Isabelle, that is, by executing the

equations stated in these theorems as term rewriting in the logic. Instead, we perform the computation outside the logic, employing an ML datatype aimed at representing efficiently the finite and the regular derivation trees dwelling the Isabelle type `dtree` from Sect. 2:

```
datatype wit_tree = Wit_Leaf of int |
                  Wit_Node of (int * int * int list) * wit_tree list
```

Here, `Wit_Node ((i, j, is), ts)` stores the root nonterminal  $i$ , a numeric identifier of the used production  $j$ , and the continuation consisting of the terminals  $is$  and the further non-terminal expanded trees  $ts$ . Moreover, `Wit_Leaf i` stores, for the case of regular infinite trees, the nonterminal where a regularity loop occurs, i.e., such that it has a previous occurrence on the path to the root.

From the ML trees, we produce witnesses represented as Isabelle constants of appropriate types (the  $w_I$ 's described above), by essentially mimicking the (co)recursive definitions employed in the proofs of the soundness parts of Theorems 4 and 5 from Sect. 4.2. Then, we certify the witnesses by producing the relevant Isabelle proof goals and discharging them by mirroring the corresponding (co)inductive arguments from the aforementioned proofs. In summary: The witness computation happens outside the logic, but the witnesses are then verified through Isabelle's kernel. After introducing any BNF, its witness set is minimized, in that redundant witnesses are removed.

The whole development pertaining to the production and certification of witnesses amounts to about 500 lines of Standard ML [8].

## 6. Related Work

Coinductive (or coalgebraic) datatypes have become popular in recent years in the study of infinite behaviors and non-terminating computation. Whereas inductive datatypes are well-studied and widely available in most programming languages and theorem provers, coinductive types are still not mainstream and pose great challenges to be integrated into current systems. Much research has appeared in the last years, in the context of theorem proving, on how to add coinductive types or improve support of coinductive proofs, including developments in Coq [7, 25], Agda [3], and CIRC [23]. The work of this paper is in line with this research and is part of the larger task of extending Isabelle/HOL with a flexible coinductive infrastructure.

Other definitional packages must also prove nonemptiness of newly defined types, but typically the proofs are easy: Homeier's quotient package for HOL4 [16] exploits the observation that quotients of nonempty sets are nonempty, and Huffman's (co)recursive domain package for Isabelle/HOLCF [18] can rely on a minimal element  $\perp$ . For the traditional datatype packages introduced by Melham [24], extended by Gunter [11], simplified by Harrison [12], and implemented in Isabelle/HOL by Berghofer and Wenzel [6], proving nonemptiness is nontrivial, but by reducing nested definitions to mutual definitions, they could employ a standard reachability analysis [6, §4.1]. To our knowledge, the completeness of the analysis has not been proved (or even formulated) for previous datatype packages.

Obviously, our overall approach to (co)datatypes is heavily inspired by category theory—this is discussed in detail in a previous paper [32], with many pointers to the literature. In category theory settings however, type nonemptiness is often not crucial, either because the developments do not target any particular logic [5, 9, 13, 15, 29] or because the targeted logics cater for (potentially empty) dependent types [1, 4, 14, 30]. Our nonemptiness-witness maintenance has similarities with the preservation of enriched types along various constructions, e.g., initial algebras and final coalgebras of pointed functors are also pointed [17]. On the other hand, none of these categorical frameworks is concerned with notions of completeness of this extra information.

## 7. Conclusion

We presented a complete solution to the nonemptiness witness problem that arises in the context of Isabelle’s new (co)datatype package. The solution exploits the functorial view of types that lies at the heart of the package. Despite its location at the intersection of category theory and higher-order logic, the problem has broad practical motivation in terms of the popularity of HOL provers. Our solution, like the rest of the definitional package, is part of the 2013 edition of Isabelle.

The problem and its solution also enjoy an elegant metatheory, which itself is best expressed in terms of codatatypes. In our case study, the package simultaneously presented itself and parts of its own metatheory.

Although the underlying categorical constructions were introduced in an earlier paper [32], we found it useful to recall them here. We used this opportunity to connect the theoretical constructions with the concrete properties made available to end users and describe the constructions in their full generality—for  $n$ -ary functors, mutual initial algebras, and mutual final coalgebras.

## Acknowledgment

Tobias Nipkow made this collaboration possible. Andreas Lochbihler suggested many improvements, notably concerning the format of concrete coinduction principles. Brian Huffman suggested major conceptual simplifications to the package. Florian Haftmann, Christian Urban, and Makarius Wenzel guided us in the jungle of package writing. Andreas Abel and Martin Hofmann pointed out relevant work. Mark Summerfield and several anonymous reviewers commented on drafts of this paper. The first author’s research was supported by the Deutsche Forschungsgemeinschaft (DFG) project Hardening the Hammer (grant Ni 491/14-1). The second author’s research was supported by the project Security Type Systems and Deduction (grant Ni 491/13-1) as part of the DFG program Reliably Secure Software Systems (RS<sup>3</sup>, Priority Program 1496). The third author’s research was supported by the DFG doctorate program 1480 (PUMA). The authors are listed in alphabetical order regardless of individual contributions or seniority.

## References

- [1] M. Abbott, T. Altenkirch, and N. Ghani. Containers: Constructing strictly positive types. *Theor. Comput. Sci.*, 342(1):3–27, 2005.
- [2] A. Abel and T. Altenkirch. A predicative strong normalisation proof for a lambda-calculus with interleaving inductive types. In T. Coquand, P. Dybjer, B. Nordström, and J. M. Smith, eds., *TYPES ’99*, vol. 1956 of *LNCS*, pp. 21–40. Springer, 2000.
- [3] A. Abel, B. Pientka, D. Thibodeau, and A. Setzer. Copatterns: programming infinite structures by observations. In *POPL*, pp. 27–38, 2013.
- [4] T. Altenkirch and P. Morris. Indexed containers. In *LICS 2009*, pp. 277–285. IEEE, 2009.
- [5] M. Barr. Terminal coalgebras in well-founded set theory. *Theor. Comput. Sci.*, 114(2):299–315, 1993.
- [6] S. Berghofer and M. Wenzel. Inductive datatypes in HOL—Lessons learned in formal-logic engineering. In Y. Bertot, G. Dowek, A. Hirschowitz, C. Paulin, and L. Théry, eds., *TPHOLs ’99*, vol. 1690 of *LNCS*, pp. 19–36, 1999.
- [7] Y. Bertot. Filters on coinductive streams, an application to Eratosthenes’ sieve. In *TLCA ’05*, pp. 102–115, 2005.
- [8] J. C. Blanchette, A. Popescu, and D. Traytel. Formal development associated with this paper. [http://www21.in.tum.de/~traytel/codata\\_wit\\_devel.tar.gz](http://www21.in.tum.de/~traytel/codata_wit_devel.tar.gz).
- [9] N. Ghani, P. Johann, and C. Fumex. Fibrational induction rules for initial algebras. In *CSL*, pp. 336–350, 2010.
- [10] E. L. Gunter. Why we can’t have SML-style datatype declarations in HOL. In L. J. M. Claesen and M. J. C. Gordon, eds., *TPHOLs ’92*, vol. A-20 of *IFIP Transactions*, pp. 561–568. North-Holland/Elsevier, 1993.
- [11] E. L. Gunter. A broader class of trees for recursive type definitions for HOL. In J. J. Joyce and C.-J. H. Seger, eds., *HUG ’93*, vol. 780 of *LNCS*, pp. 141–154. Springer, 1994.
- [12] J. Harrison. Inductive definitions: Automation and application. In E. T. Schubert, P. J. Windley, and J. Alves-Foss, eds., *TPHOLs ’95*, vol. 971 of *LNCS*, pp. 200–213. Springer, 1995.
- [13] R. Hasegawa. Two applications of analytic functors. *Theor. Comput. Sci.*, 272(1-2):113–175, 2002.
- [14] U. Hensel and B. Jacobs. Proof principles for datatypes with iterated recursion. In *Category Theory and Computer Science*, pp. 220–241, 1997.
- [15] C. Hermida and B. Jacobs. Structural induction and coinduction in a fibrational setting. *Inf. Comput.*, 145(2):107–152, 1998.
- [16] P. V. Homeier. A design structure for higher order quotients. In J. Hurd and T. F. Melham, eds., *TPHOLs 2005*, vol. 3603 of *LNCS*, pp. 130–146. Springer, 2005.
- [17] B. T. Howard. Inductive, coinductive, and pointed types. In *ICFP*, pp. 102–109, 1996.
- [18] B. Huffman. A purely definitional universal domain. In S. Berghofer, T. Nipkow, C. Urban, and M. Wenzel, eds., *TPHOLs 2009*, vol. 5674 of *LNCS*, pp. 260–275. Springer, 2009.
- [19] The Isabelle theorem prover, 2013. <http://isabelle.in.tum.de/>.
- [20] A. Krauss and T. Nipkow. Proof pearl: Regular expression equivalence and relation algebra. *J. Autom. Reasoning*, 49(1):95–106, 2012.
- [21] M. Lenisa, J. Power, and H. Watanabe. Distributivity for endofunctors, pointed and co-pointed endofunctors, monads and comonads. *Electr. Notes Theor. Comput. Sci.*, 33:230–260, 2000.
- [22] A. Lochbihler. Java and the Java memory model—A unified, machine-checked formalisation. In H. Seidl, ed., *ESOP 2012*, vol. 7211 of *LNCS*, pp. 497–517. Springer, 2012.
- [23] D. Lucanu, E.-I. Goriac, G. Caltais, and G. Roşu. CIRC: A behavioral verification tool based on circular coinduction. In *CALCO’09*, pp. 433–442, 2009.
- [24] T. F. Melham. Automating recursive type definitions in higher order logic. In G. Birtwistle and P. A. Subrahmanyam, eds., *Current Trends in Hardware Verification and Automated Theorem Proving*, pp. 341–386. Springer, 1989.
- [25] K. Nakata, T. Uustalu, and M. Bezem. A proof pearl with the fan theorem and bar induction—Walking through infinite trees with mixed induction and coinduction. In H. Yang, ed., *APLAS ’11*, vol. 7078 of *LNCS*, pp. 353–368. Springer, 2011.
- [26] L. C. Paulson. A fixedpoint approach to (co)inductive and (co)datatype definitions. In G. D. Plotkin, C. Stirling, and M. Tofte, eds., *Proof, Language, and Interaction—Essays in Honour of Robin Milner*, pp. 187–212. MIT Press, 2000.
- [27] B. C. Pierce. *Types and Programming Languages*. MIT Press, 2002.
- [28] J. J. M. M. Rutten. Relators and metric bisimulations. *Electr. Notes Theor. Comput. Sci.*, 11:252–258, 1998.
- [29] J. J. M. M. Rutten. Universal coalgebra: A theory of systems. *Theor. Comput. Sci.*, 249:3–80, 2000.
- [30] C. Schwaab and J. G. Siek. Modular type-safety proofs in Agda. To appear in *PLPV 2013*.
- [31] M. B. Smyth. Power domains. *J. Comput. Syst. Sci.*, 16(1):23–36, 1978.
- [32] D. Traytel, A. Popescu, and J. C. Blanchette. Foundational, compositional (co)datatypes for higher-order logic—Category theory applied to theorem proving. In *LICS 2012*, pp. 596–605. IEEE, 2012.

## APPENDIX

Here we give more details concerning the concepts discussed in the paper, including proofs.

### A. Embedding of Initial Algebras into Final Coalgebras

It is well-known that, for  $n$  ( $m+n$ ) functors  $(\bar{\alpha}, \bar{\beta}) F_j$ , their mutual initial algebra can be embedded in their mutual final coalgebra via  $\text{emb}_j : \bar{\alpha} \text{IF}_j \rightarrow \bar{\alpha} \text{JF}_j$  defined

- either recursively, as  $\text{emb}_j = \text{fold}_j \overline{\text{ctor}'}$ , where  $\text{ctor}'_j : (\bar{\alpha}, \bar{\alpha} \text{JF}) F_j \rightarrow \bar{\alpha} \text{JF}_j$  is the inverse of the final-coalgebra unfolding bijection  $\text{ctor}_j$ ,
- or corecursively, as  $\text{emb}_j = \text{unfold}_j \overline{\text{ctor}'}$ , where  $\text{ctor}'_j : \bar{\alpha} \text{IF}_j \rightarrow (\bar{\alpha}, \bar{\alpha} \text{IF}) F_j$  is the inverse of the initial-algebra folding bijection  $\text{ctor}_j$ .

Along this embedding,  $\text{ctor}'_j$  is the extension of  $\text{ctor}_j$  and  $\text{ctor}_j$  is the extension of  $\text{ctor}'_j$ . Moreover, the image of  $\text{emb}_j$  consists precisely of the “finite” elements of  $\bar{\alpha} \text{JF}_j$ , i.e., of those characterized by the induction principle of  $\bar{\alpha} \text{IF}_j$ . Therefore, we define  $\overline{\text{JFfinite}}$  as the least tuple of predicates satisfying the following specification.

$$\forall x \in (\bar{\alpha}, \bar{\alpha} \text{IF}) F_j. (\bigwedge_{k=1}^n \forall b \in \text{Fset}_j^{m+k} x. \text{JFfinite}_k b) \Rightarrow \text{JFfinite}_j(\text{ctor}_j x)$$

We have that, for all  $j$ ,  $\text{emb}_j \cdot (\bar{\alpha} \text{IF}_j) = \{a \in \bar{\alpha} \text{JF}_j \mid \text{JFfinite}_j a\}$ .

### B. More on the Coinductive Tree Case Study

Next we define the notions of interior and frontier of a tree.  $\text{ltr} : \text{dtree} \rightarrow \text{N set}$  is defined inductively as follows:

$$\begin{aligned} \text{root } t &\in \text{ltr } t \\ \text{Inr } t_1 \in \text{cont } t \wedge n \in \text{ltr } t_1 &\Rightarrow n \in \text{ltr } t \end{aligned}$$

$\text{Fr} : \text{dtree} \rightarrow \text{N set}$  is defined inductively as follows:

$$\begin{aligned} \text{Inl } t \in \text{cont } t &\Rightarrow t \in \text{Fr } t \\ \text{Inr } t_1 \in \text{cont } t \wedge t \in \text{Fr } t_1 &\Rightarrow t \in \text{Fr } t \end{aligned}$$

The parameterized versions of interior and frontier (as required in Sect. 4.3),  $\text{ltr} : \text{N set} \rightarrow \text{dtree} \rightarrow \text{N set}$  and  $\text{Fr} : \text{N set} \rightarrow \text{dtree} \rightarrow \text{N set}$ , are defined inductively modifying the above clauses textually as follows:

- ‘ltr’ is replaced by ‘ltr ns’ and ‘Fr’ is replaced by ‘Fr ns’;
- the hypothesis  $\text{root } t \in \text{ns}$  is added.

The notion of a finite tree,  $\text{ftree} : \text{dtree} \rightarrow \text{bool}$ , is also defined inductively:

$$\begin{aligned} \text{Inr} \cdot \text{cont } t = \emptyset &\Rightarrow \text{ftree } t \\ (\forall t_1 \in \text{Inr}^- t. \text{ftree } t_1) &\Rightarrow \text{ftree } t \end{aligned}$$

$\text{ftree}$  is in effect an instance of the  $\text{JFfinite}$  predicate introduced in Appendix A.

### C. Inductive Trees

The datatype definition

$$\text{datatype } \text{fdtree} = \text{FNode} (\text{froot} : \text{N}) (\text{fcont} : (\text{T} + \text{dtree}) \text{fset})$$

(introducing finite trees) produces the operations  $\text{FNode}$ ,  $\text{froot}$ , and  $\text{fcont}$  having constructor and selector properties corresponding precisely to the ones of  $\text{Node}$ ,  $\text{root}$  and  $\text{cont}$  from the codatatype  $\text{dtree}$  in Sect. 2.1.

The difference concerns induction and recursion.

### Iteration

The general principle described in Sect. 3.1 employs in this unary case a iterator fold of (polymorphic) type  $(\beta \text{ pre\_fdtree} \rightarrow \beta) \rightarrow \text{fdtree} \rightarrow \beta$ , for which it yields

$$\forall s : \beta \text{ pre\_fdtree} \rightarrow \beta. (\text{fold } s) \circ \text{ctor} = s \circ \text{pre\_fdtree\_map} (\text{fold } s)$$

that is,

$$\forall s : \beta \text{ pre\_fdtree} \rightarrow \beta. \forall k. \text{fold } s (\text{ctor } k) = s (\text{pre\_fdtree\_map} (\text{fold } s) k)$$

The  $\text{fdtree}$ -defining BNF is the same as the  $\text{dtree}$ -defining BNF:  $\beta \text{ pre\_fdtree} = \text{N} \times (\text{T} + \beta) \text{fset}$  and  $\text{pre\_fdtree\_map } f = \text{id} \otimes (\text{image} (\text{id} \oplus f))$ .

As in the codatatype case, the above characterization needs some customization. Using the  $\text{FNode}$  instead of  $\text{ctor}$  and unfolding the definition of  $\text{pre\_fdtree\_map}$ , we obtain

$$\begin{aligned} \forall s : \text{N} \times (\text{T} + \beta) \text{fset} \rightarrow \beta. \forall n \text{ as}. \\ \text{fold } s (\text{FNode } n \text{ as}) = s (\text{pre\_fdtree\_map} (\text{fold } s) (n, \text{as})) \end{aligned}$$

By unfolding the definition of  $\text{pre\_fdtree\_map}$ , we obtain

$$\begin{aligned} \forall s : \text{N} \times (\text{T} + \beta) \text{fset} \rightarrow \beta. \forall n \text{ as}. \\ \text{fold } s (\text{FNode } n \text{ as}) = s (n, (\text{id} \oplus \text{fold } s) \cdot \text{as}) \end{aligned}$$

Finally, replacing  $\text{fold}$  with its more convenient curried version  $\text{fold}' : (\text{N} \rightarrow (\text{T} + \beta) \text{fset} \rightarrow \beta) \rightarrow \text{fdtree} \rightarrow \beta$  defined as  $\text{fold}' s = \text{fold} (\lambda(n, \text{as}). s \text{ n as})$ , we obtain the following customized iteration principle, where we write  $\text{fold}$  instead of  $\text{fold}'$ :

For all sets  $\beta$ , functions  $s : \text{N} \rightarrow (\text{T} + \beta) \text{fset} \rightarrow \beta$  and elements  $n \in \text{N}$  and  $\text{as} \in (\text{T} + \text{fdtree}) \text{fset}$ , it holds that  $\text{fold } s (\text{FNode } n \text{ as}) = s \text{ n } ((\text{id} \oplus \text{fold } s) \cdot \text{as})$ .

### Induction

The induction principle from Sect. 3.3 yields for  $\varphi : \alpha \text{fdtree} \rightarrow \text{bool}$

$$\frac{\forall k \in \alpha \text{ pre\_fdtree}. (\forall t \in \text{Fset } k. \varphi t) \Rightarrow \varphi (\text{ctor } k)}{\forall t \in \alpha \text{fdtree}. \varphi t}$$

i.e., using the curried variation  $\text{FNode}$  of  $\text{ctor}$ ,

$$\frac{\forall n \text{ as}. (\forall t \in \text{Fset} (n, \text{as}). \varphi t) \Rightarrow \varphi (\text{FNode } n \text{ as})}{\forall t \in \alpha \text{fdtree}. \varphi t}$$

Unfolding the definition of  $\text{Fset}$ , namely,  $\text{Fset} (n, \text{as}) = \text{Inr}^- \text{as}$ , we obtain the end-product customized induction for finite trees:

$$\frac{\forall n \text{ as}. (\forall t \in \text{Inr}^- \text{as}. \varphi t) \Rightarrow \varphi (\text{FNode } n \text{ as})}{\forall t \in \alpha \text{fdtree}. \varphi t}$$

### Embedding of fdtree in dtree

Here is what the standard embedding  $\text{emb}$  described in Appendix A looks like for the concrete instance  $\text{fdtree} \rightarrow \text{dtree}$ :

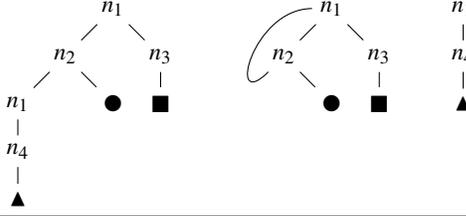
- recursively:  $\text{fold } \text{Node}$ , yielding  $\text{emb} (\text{FNode } n \text{ as}) = \text{Node } n ((\text{id} \oplus \text{emb}) \cdot \text{as})$ ;
- corecursively:  $\text{unfold } \text{froot } \text{fcont}$ , yielding the equations  $\text{root} (\text{emb } t) = \text{froot } t$  and  $\text{cont} (\text{emb } t) = (\text{id} \oplus \text{emb}) \cdot \text{fcont } t$ .

As an instance of the situation from Appendix A, the image of  $\text{emb}$  is the set of finite  $\text{dtree}$ 's,  $\{t \in \text{dtree}. \text{ftree}\}$ .

### D. Proofs

For more details on some of the proofs, we refer the reader to our Isabelle formalization [8], which employs essentially the same notations as this text.

*Proof of Theorem 1:*  $\text{H } n$  is regular by construction: if a subtree of it has root  $n'$ , then it is equal to  $\text{H } n'$ . The frontier inclusion



**Fig. 3.** A finite derivation tree (left), a regular cut of it (middle), and a finite regular cut of it (right)

$\text{Fr}(\text{H}n) \subseteq \text{Fr}t_0$  follows by routine fixpoint induction on the definition of  $\text{Fr}$  (since at each node  $n' \in \text{ltr}(\text{H}n)$  we only have the immediate leaves of  $\text{pick } n'$ , which is a subtree of  $\text{Fr}t_0$ ). Finally, assume that  $t_0$  is well-formed. Then the fact that  $\text{H}n$  is well-formed follows by routine fixpoint coinduction on the definition of  $\text{wf}$  (since, again, at each  $n' \in \text{ltr}(\text{H}n)$  we have the production of  $\text{pick } n'$ ).  $\square$

*Proof of Theorem 2:* Let  $\mathcal{K} = \{\bigcup_{j \in J} I_j \mid J \in \mathcal{J} \wedge (I_j)_j \in \prod_{j \in J} \mathcal{S}_j\}$ . We first prove that  $\mathcal{K}$  is sound for  $\text{H}$ . Let  $K \in \mathcal{K}$  and  $\bar{\alpha}_m$  such that  $\forall i \in K. \alpha_i \neq \emptyset$ . By the definition of  $\mathcal{K}$ , we obtain  $J \in \mathcal{J}$  and  $(I_j)_{j \in J}$  such that (1)  $K = \bigcup_{j \in J} I_j$  and (2)  $\forall j \in J. I_j \in \mathcal{S}_j$ . Using (1), we have  $\forall j \in J. \forall i \in I_j. \alpha_i \neq \emptyset$ . Hence, since each  $\mathcal{S}_j$  is sound for  $F_j$ ,  $\forall j \in J. \bar{\alpha} F_j \neq \emptyset$ . Finally, since  $\mathcal{J}$  is sound for  $\text{G}$ , we obtain  $\bar{\alpha} \bar{F} \text{G} \neq \emptyset$ , i.e.,  $\bar{\alpha} \text{H} \neq \emptyset$ , as desired.

We now prove that  $\mathcal{K}$  is complete for  $\text{H}$ . Let  $K \subseteq [m]$  be a  $\text{H}$ -witnessed set of indices. Let  $\bar{\beta}_n$  be defined as  $\beta_j = \text{unit}$  if  $j \in K$  and  $= \emptyset$  otherwise and let  $J = \{j \in [n]. \bar{\beta} F_j \neq \emptyset\}$ . Since  $K$  is  $\text{H}$ -witnessed, we obtain that  $\bar{\beta} \text{H} \neq \emptyset$ , i.e., (1)  $\bar{\beta} \bar{F} \text{G} \neq \emptyset$ .

We show that (3)  $\text{G}$  is  $J$ -witnessed. Let  $\bar{\gamma}_n$  such that  $\forall j \in J. \gamma_j \neq \emptyset$ . Thanks to the definition of  $J$ , we have  $\forall j \in [n]. F_j \neq \emptyset \Rightarrow \gamma_j \neq \emptyset$ , and therefore we obtain the functions  $(f_j : \bar{\beta} F_j \rightarrow \gamma_j)_{j \in [n]}$ . Then, thanks to  $\text{Gmap } \bar{f} : \bar{\beta} \bar{F} \text{G} \rightarrow \bar{\gamma} \text{G}$ , by (1) we obtain that  $\bar{\gamma} \text{G} \neq \emptyset$ , as desired.

From (3), since  $\mathcal{J}$  is complete for  $J$ , we obtain  $J_1 \in \mathcal{J}$  such that  $J_1 \subseteq J$ . Let  $j \in J_1$ . By the definition of  $J$ , we have  $\bar{\beta} F_j \neq \emptyset$ , making  $\bar{\beta} F_j$   $K$ -witnessed (by the definition of  $\bar{\beta}$ ); hence, since  $\mathcal{S}_j$  is  $F_j$ -complete, we obtain  $I_j \in \mathcal{S}_j$  such that  $I_j \subseteq K$ . Then  $K_1 = \bigcup_{j \in J_1} I_j$  belongs to  $\mathcal{K}$  and is included in  $K$ , as desired.  $\square$

*Proof of Theorem 3:* Immediate from the definitions.  $\square$

In some of the following proofs we exploit an embedding of datatypes as finite codatatypes (Appendix A). Using this embedding, we can transfer the recursive definition and structural induction principles from  $\bar{\text{IF}}$  to finite elements of  $\bar{\text{JF}}$ , and in particular from  $\text{fdtree}$  to finite trees in  $\text{dtree}$ .

The regular cut of a tree works well with respect to the codatatype metatheory, but for datatypes it has the disadvantage that it may produce infinite trees out of finite ones (cf. Figure 3, left and middle). We need a slightly different concept for datatypes: the finite regular cut. Let  $t_0$  be a finite derivation tree. We choose the function  $\text{fpick} : \text{ltr } t_0 \rightarrow \text{Subtr } t_0$  similarly to  $\text{pick}$  from Sect. 2.3, but making sure that in addition the choice of the subtrees  $\text{fpick } n$  is minimal, in that  $\text{fpick } n$  does not have  $n$  in the interior of a proper subtree (and hence does not have any proper subtree of root  $n$ )—such a choice is possible thanks to the finiteness of  $t_0$ . We define the finite regular cut of  $t_0$ ,  $\text{rfcut } t_0$ , just like  $\text{rcut } t_0$  but using  $\text{fpick}$  instead of  $\text{pick}$ . Now we can prove:

**Lemma 1.** Assume  $t_0$  is a finite derivation tree. Then:

- (1) The statement of Theorem 1 holds if we replace  $\text{rcut}$  by  $\text{rfcut}$ .
- (2)  $\text{rfcut } t_0$  is finite.

*Proof:* (1) Similar to the proof of Theorem 1. (2) By routine induction on  $t_0$ .  $\square$

*Proof of Theorem 5:* Let  $j_0 \in [n]$ . We first show that  $\mathcal{L}_{\text{Gr}}^{\text{rf}}(j_0)$  is sound. Let  $t_0$  be a well-formed finite regular derivation tree with root  $j_0$ . We need to prove that  $F_{j_0}$  is  $\text{Fr } t_0$ -witnessed. For this, we fix  $\bar{\alpha}_m$  such that  $\forall i \in \text{Fr } t_0. \alpha_i \neq \emptyset$ , and aim to show that  $\bar{\alpha} \text{IF}_{j_0} \neq \emptyset$ .

For each  $j \in \text{ltr } t_0$ , let  $t_j$  be the corresponding subtree of  $t_0$ . (It is well-defined, since  $t_0$  is regular.) Note that  $t_0 = t_{j_0}$ . For each  $K$  such that  $(j, \text{cp}(K)) \in \text{P}$ , since  $K \in \mathcal{K}_j$  and  $\mathcal{K}_j$  is sound for  $F_j$ , we obtain a  $K$ -witness for  $F_j$ , i.e., a function  $w_{j,K} : (\alpha_k)_{k \in K} \rightarrow \bar{\alpha} F_j$ .

We verify the following fact by induction on the finite derivation tree  $t$ : If  $\exists j \in \text{ltr } t_0. t = t_j$ , then  $\bar{\alpha} \text{IF}_j \neq \emptyset$ . The induction step goes as follows: Assume  $t = t_j$  has the form  $\text{Node } j \text{ as}$ , and let  $J$  be the set of all roots of the immediate subtrees of  $t$ , namely,  $\text{root} \cdot (\text{Inr}^- \text{cont } t)$ . By the induction hypothesis,  $\bar{\alpha} \text{IF}_{j'} \neq \emptyset$  (say,  $b_{j'} \in \bar{\alpha} \text{IF}_{j'}$ ) for all  $j' \in J$ . Then  $w_{j,K}(\alpha_i)_{i \in \text{Inr}^- t} (b_{j'})_{j' \in J} \in \bar{\alpha} \text{IF}_j$ , making  $\bar{\alpha} \text{IF}_j$  nonempty.

In particular,  $\bar{\alpha} \text{JF}_{j_0} \neq \emptyset$ , as desired.

We now show that  $\mathcal{L}_{\text{Gr}}^{\text{rf}}(j_0)$  is complete. Let  $I \subseteq [m]$  such that  $\text{IF}_{j_0}$  is  $I$ -witnessed. We need to find  $I_1 \in \mathcal{L}_{\text{Gr}}^{\text{rf}}(j_0)$  such that  $I_1 \subseteq I$ .

Let  $\bar{\alpha}_m$  be defined as  $\alpha_i = \text{unit}$  if  $i \in I$  and  $= \emptyset$  otherwise. We verify, by structural  $\bar{\text{IF}}$ -induction on  $b$ , that for all  $j \in [n]$  and  $b \in \bar{\alpha} \text{IF}_j$ , there exists a finite well-formed derivation tree  $t$  such that  $\text{root } t = j$  and  $\text{Fr } t \subseteq I$ . For the inductive step, assume  $\text{ctor } j \ x \in \bar{\alpha} \text{IF}_j$ , where  $x \in (\bar{\alpha}, \bar{\alpha} \bar{\text{IF}}) F_j$ . By the induction hypotheses, we obtain the finite well-formed derivation trees  $\bar{t}_n$  such that  $\text{root } \bar{t}_j = j$  and  $\text{Fr } \bar{t}_j \subseteq I$  for all  $j \in [n]$ . Let  $J = \{j' \in [n]. \bar{\alpha} \text{IF}_{j'} \neq \emptyset\}$ . Then  $F_j$  is  $(I \cup J)$ -witnessed, hence by the  $F_j$ -completeness of  $\mathcal{K}_j$  we obtain  $K \in \mathcal{K}_j$  such that  $K \subseteq I \cup \{m + j'\}. j' \in J\}$ . We take  $t$  to have  $j$  as root,  $I \cap K$  as leaves and  $(\bar{t}_{j'})_{j' \in J}$  as immediate subtrees; namely,  $t = \text{Node } j \ ((\text{Inl} \cdot I) \cup (\text{Inr} \cdot \{t_{j'} \cdot j' \in J\}))$ .

Let  $t_0$  be a tree as above corresponding to  $j_0$  (since  $\bar{\alpha} \text{IF}_{j_0} \neq \emptyset$ ). Then, by Lemma 1,  $t_1 = \text{rcut } t_0$  is a well-formed finite derivation tree such that  $\text{Fr } t_1 \subseteq \text{Fr } t_0 \subseteq I$ . Thus, taking  $I_1 = \text{Fr } t_1$ , we obtain  $I_1 \in \mathcal{L}_{\text{Gr}}^{\text{rf}}(j_0)$  and  $I_1 \subseteq I$ , as desired.  $\square$

*Proof of Theorem 7:*  $\mathcal{L}_{\text{Gr}}^{\text{r}} \text{ns } n \subseteq \{\emptyset\}$ , since  $\text{Fr } \text{ns } t = \emptyset$  for all  $t$  such that  $\text{root } t = n$ . It remains to show  $\emptyset \in \mathcal{L}_{\text{Gr}}^{\text{r}} \text{ns } t$ , i.e., to find a derivation tree with root  $n$ . In fact, using the assumption that there are no unused nonterminals, we can build a “default derivation tree”  $\text{deftr } n$  for each  $n$  as follows. We pick, for each  $n$ , a set  $S \ n \in (\text{T} + \text{N}) \text{fset}$  such that  $(n, S \ n) \in \text{P}$ . Then we define  $\text{deftr} : \text{N} \rightarrow \text{dtree}$  corecursively as  $\text{deftr} = \text{unfold id } S$ , i.e., such that  $\text{root}(\text{deftr } n) = n$  and  $\text{cont}(\text{deftr } n) = (\text{id} \oplus \text{deftr}) \cdot S \ n$ . It is easy to prove by  $\text{KT}$ -coinduction that  $\text{deftr } n$  is a derivation tree for each  $n$ .

Now assume  $n \notin \text{ns}$ , and let  $\text{ns}' = \text{ns} \setminus \{n\}$ .

For the left-to-right direction, we prove more than  $\leq$ , namely, actual inclusion between  $\mathcal{L}_{\text{Gr}}^{\text{r}} \text{ns } n$  and the righthand side. Assume  $t$  is a well-formed regular derivation tree of root  $n$ . We need to find  $ss \in (\text{T} + \text{N}) \text{fset}$  and  $U : \text{Inr}^- ss \rightarrow \text{dtree}$  such that, for all  $n' \in \text{Inr}^- ss$ ,  $U \ n'$  is a well-formed regular derivation tree of root  $n'$  and  $\text{Fr } \text{ns } t = \text{Inl}^- ss \cup \bigcup_{n' \in \text{Inr}^- ss} \text{Fr } \text{ns}'(U \ n')$ .

Clearly  $ss$  should be the right-hand side of the top production of  $t$ . As for  $U$ , of course the immediate subtrees of  $t$  provide intuitive candidates; however, these do not work, since our goal is to have  $\text{Fr } \text{ns } t$  covered by  $(\text{Inl}^- ss$  in conjunction with)  $\text{Fr } \text{ns}'(U \ n')$ , while the immediate subtrees only guarantee this property with respect to  $\text{Fr } \text{ns}'(U \ n')$ , i.e., allowing paths to go through  $n$  as well. A correct solution is again offered by a corecursive definition: We build the tree  $t_0$  from  $t$  by substituting hereditarily each subtree with root  $n$  by  $t$ . Formally, we take  $t_0 = \text{unfold } r \ c$ , where  $r \ t' = \text{root } t'$  and  $c \ t' = \text{cont } t'$  if  $\text{root } t' = n$  and  $c \ t' = \text{cont } t'$  otherwise. It is easy to prove that  $t_0$ , like  $t$ , is a regular derivation tree. Thus, we can define  $U$  to give, for any  $n'$ , the corresponding immediate subtree of  $t_0$ .

To prove the right-to-left direction, let  $ss \in (\mathbb{T} + \mathbb{N})$  fset and  $K \in \prod_{n' \in \text{Inr}^- ss} \mathcal{L}_{\text{Gr}}^r ns' n'$  such that  $ts = \text{Inl}^- ss \cup \bigcup_{n' \in \text{Inr}^- ss} K_{n'}$ . Unfolding the definition of  $\mathcal{L}_{\text{Gr}}^r$ , we obtain  $U : \text{Inr}^- ss \rightarrow \text{dtree}$  such that, for all  $n' \in \text{Inr}^- ss$ ,  $U n'$  is a regular derivation tree of root  $n'$  such that  $K_{n'} \in \text{Fr} ns' (U n')$ . Then the tree of immediate leafs  $\text{Inl}^- ss$  and immediate subtrees  $\{U n'. n' \in \text{Inr}^- ss\}$ , namely,  $\text{Node } n ((\text{id} \oplus U) \cdot ss)$ , is the desired regular derivation tree whose frontier is included  $ts$ .  $\square$

In what follows,  $nl$  ranges over lists of nonterminals and ‘ $\cdot$ ’ denotes list concatenation. If  $n$  is a nonterminal,  $n$  also denotes the  $n$ -singleton list.

The predicate  $\text{path } nl \ t$ , stating that  $nl$  is a path in  $t$  (starting from the root), is defined inductively by the following rules:

$$\begin{array}{l} \text{path } (\text{root } t) \ t \\ \text{Inr } t' \in \text{cont } t \wedge \text{path } nl \ t' \Rightarrow \text{path } ((\text{root } t) \cdot nl) \ t' \end{array}$$

**Lemma 2.** Let  $t$  be a finite regular derivation tree. Then  $t$  has no paths that contain repetitions.

*Proof:* Assume, by absurdity, that a path  $nl$  in  $t$  contains repetitions, i.e., has the form  $nl_1 \cdot n \cdot nl_2 \cdot n$ , and let  $t_1$  and  $t_2$  be the subtrees corresponding to the paths  $nl_1 \cdot n$  and  $nl_2$ , respectively. Then  $t_2$  is a proper subtree of  $t_1$ ; on the other hand, by the regularity of  $t$ , we have  $t_1 = t_2$ , which is impossible since  $t_1$  and  $t_2$  are finite.  $\square$

*Proof of Theorem 8:* According to Lemma 2 and the properties of regular cuts, we have that (1)  $\mathcal{L}_{\text{Gr}}^r ns' n \equiv \mathcal{L}_{\text{Gr}}^{\text{pf}} ns' n$ , where  $\mathcal{L}_{\text{Gr}}^{\text{pf}} ns' n$  is the language defined just like  $\mathcal{L}_{\text{Gr}}^r ns' n$ , but replacing “regular” with “having no paths that contain repetitions.” Moreover, it is easy to see that (2) the desired facts hold if we replace  $\mathcal{L}_{\text{Gr}}^r ns' n$  with  $\mathcal{L}_{\text{Gr}}^{\text{pf}} ns' n$  and  $\equiv$  with equality. From (1) and (2) the result follows.  $\square$

## E. Registration of a New BNF in Isabelle

The type constructor  $\alpha$  bag of bags (multisets) is registered through the following command:

```
bnf_def
  bag_map :: ('a => 'b) => 'a bag => 'b bag
  [set_of :: 'a bag => 'a set]
  natLeq :: (nat * nat) set
  [{#} :: 'a bag]
  bag_rel :: ('a => 'b => bool) =>
    'a bag => 'b bag => bool
```

The command provides the necessary infrastructure that makes  $\alpha$  bag a BNF, consisting of various previously introduced constants (whose definition is not shown here):

- the functorial action, `bag_map`;
- the “Fset” operation discussed in Sect. 3.3, `set_of`;
- a cardinal bound, here, that of natural numbers, `natLeq` (in Isabelle, cardinals are represented as minimal well-order relations);
- a witness term, here, the empty bag `{#}`;
- a customized relator, as discussed in Sect. 3.4, `bag_rel`.

Then the user is prompted to prove a few facts, including one saying that `set_of` is a natural transformation:

$$\text{set\_of } o \ \text{bag\_map } f = \text{image } f \ o \ \text{set\_of}$$

Upon discharging these goals, the  $\alpha$  bag BNF is registered.