

# A Recursive Inclusion Checker for Recursively Defined Subtypes

Hans de Nivelle  
School of Engineering and Digital Sciences,  
Nazarbayev University, Nur-Sultan City, Kazakhstan,  
hans.denivelle@nu.edu.kz

5 декабря 2021 г.

## Аннотация

We present an inclusion checker for recursively defined subtypes of recursively defined types and apply this inclusion checker when resolving ambiguous names. This work is part of a project to design a programming language that is suitable for implementation of logic. Logical formulas are tree-like structures with many constructors having different arities and argument types. Algorithms that use these structures must perform case analysis on the constructors, and access subtrees whose type and existence depend on the constructor used.

In many programming languages, case analysis is handled by matching, but we want to take a different approach, based on subtypes. Instead of matching a tree into different possible forms, we want to classify it against a set of disjunct subtypes. Subtypes are more general than structural forms, and in addition, can be used independently as preconditions or postconditions of procedures.

In order to make this possible, we need to be able to evaluate subtypes on concrete data, so we will define a precise semantics. For the resolution of overloads between ambiguous field references, we need to be able to efficiently check propositional relations between subtypes. We define the types and the subtypes (which we will call *adjectives*), define their semantics, and give a tableaux-based inclusion checker for adjectives. We show how to use this inclusion checker for resolving ambiguous field references in declarations of adjectives.

## 1. Introduction

Our goal is to develop a programming language that is convenient for the implementation of logic. Traditionally, functional languages are considered the most suitable for implementation of logic, like OCaml ([1]), Haskell ([2]) or Scala ([3]). Functional languages have inductive types, and use matching for accessing subtrees. Matching can be viewed as simultaneously inspecting a tree and extracting subtrees into local variables. In our proposed language, we want to replace this mechanism by usual field access, as used in imperative languages like Java or C++. In order to classify a formula, so that we know which fields exist, we use subtypes. The advantage of using subtypes is that, once one has the mechanism for defining and verifying them, they are more general than structure matching, and can be used at many other points in the programming language. For example, subtypes can be used as preconditions or postconditions of functions.

As an example of our approach, suppose one has a logical formula  $f$ , and one wants to know if it constructed by **and**. One can either match it into  $\mathbf{and}(F_1, F_2)$ , or define a subtype 'constructed by **and**', and define two fields  $f_1, f_2$  that can be accessed only on formulas refined by subtype 'constructed by **and**'. The advantage of subtypes is that they can be easily combined by Boolean operators. For example, it is easy to define the subtype 'formed by **and, or**', or 'is a literal' and use this as a condition in case analysis.

We will define a system for defining recursive types, together with subtypes. We call the subtypes *adjectives*. In addition to structural matchings of fixed depth, our adjective system also allows recursively

defined subtypes. This makes it possible for example to define negation normal form or conjunctive normal form as a subtype of formulas and use it in case analysis, or as precondition or postcondition of a function.

Adjectives are somewhat similar to liquid types ([4]) or refinement types ([5]) but they have a different aim. Adjectives are not intended for checking operations whose safety depends on arithmetic. Because of this, adjectives do not rely on SMT solving. We prefer the term *adjective* over *subtype* or *refinement type*, because adjectives can overlap, and are frequently created for one time use. We also want to use adjectives for static type checking and overload resolution. In order to show that this is possible, we define a procedure that resolves ambiguous field references in the adjectives themselves.

We start by defining the type and adjective system in this section. In Section 2 we give the semantics, which can be used for evaluating adjectives on data, so that they can be used as a replacement for matching. In Section 3, we give a tableaux-based procedure for deciding satisfiability of adjectives. This procedure can be used for checking exhaustiveness of a given case analysis, for checking exclusiveness of cases, and for static type checking. In Section 4 we will apply the tableaux procedure on the problem of resolving ambiguous field references in the adjective declarations themselves. A similar algorithm could be used for resolving ambiguous function calls. We start by defining primitive types:

**Definition 1.1.** *The primitive types are bool, char, nat, double, and selector. We assume that all primitive types, with the exception of selector, have an order  $<$  defined on them.*

Type selector consists of named constants. It plays a role comparable to constructors in functional languages. In Rust ([6]), enumeration types play the same role. We write elements of selector as identifiers preceded by a question mark, e.g. ?and, ?or, ?implies, ?equiv, or ?zero, ?succ. We give an example of their use below, after the definition of our type system.

We first define types, and after that adjectives. Strictly seen, the definition of adjectives should come first, because types contain adjectives while adjectives do not contain types. We find, however, that this order would be unintuitive, because it would be hard to understand the usefulness of adjectives, without having seen the definition of types. Because of this, we start by introducing types first:

**Definition 1.2.** *We first define simple types:*

- *A primitive type is a simple type.*
- *An identifier is a simple type (assuming that it has been defined as type, see Definition 1.3 below).*
- *If  $T$  is a simple type,  $A$  is an adjective, then  $T \circ A$  is also a simple type.*

*Next come compound types. They are recursively defined as follows:*

1. *If  $v_1, \dots, v_n$  are identifiers,  $V_1, \dots, V_n$  are simple types with  $n \geq 0$ , then  $(v_1 : V_1, \dots, v_n : V_n)^*$  is a compound type.*
2. *If  $C$  is a compound type,  $v$  is an identifier and  $V$  is a simple type, then  $v : V, C$  is a compound type as well.*
3. *If  $C_1, \dots, C_m$  with  $m \geq 1$  are compound types,  $A_1, \dots, A_m$  are adjectives,  $s$  is an identifier,  $S$  is a simple type, then  $s?(A_1 \Rightarrow C_1, \dots, A_m \Rightarrow C_m)$  is a compound type as well.*

*In case 1, the fields  $v_1, \dots, v_n$  are called repeated fields. All other fields are called scalar fields. In case 3, field  $s$  is called the pivot field. Case 3 defines compound types where the existence of later fields depends on the adjective that applies to  $s$ .*

We distinguish between simple and compound types, because we want every compound type to have a name, which solves some technical problems. It has no consequences for expressiveness, because one can always define a compound type, and use the name as simple type. This can be done by the compiler, if needed.

The meaning of  $T \circ A$  is: Type  $T$  refined by adjective  $A$ . Compound types have a tree structure where every path either ends in  $()^*$  (no repeated fields) or in  $(v_1 : V_1, \dots, v_n : V_n)^*$  with  $n > 0$  (ends in  $n$  repeated

fields). Repeated fields implicitly define an array whose members must be accessed with array subscript notation.

For the moment, we assume that different fields of different types have distinct names. This is an unrealistic restriction for a real programming language, because field names can be reused in different types. In Section 4 we will allow reuse of field names between different options of a type, and between different types. We will give a procedure that is able to resolve such ambiguous uses of a field name into unambiguous names.

Type definitions need to be stored in a mapping. It is convenient to use two distinct mappings, one for the simple types, and one for the compound types:

**Definition 1.3.** We define two mappings  $\Sigma_S$  and  $\Sigma_C$ . The first mapping maps identifiers to simple types, and the second mapping maps identifiers to compound types. If  $\Sigma_S$  contains a value for  $v$ , we write  $\Sigma_S(v)$  for the value. Similarly, if  $\Sigma_C$  contains a value for  $v$ , we write  $\Sigma_C(v)$  for the value. We assume that  $\Sigma_S$  and  $\Sigma_C$  do not contain a value for the same variable  $v$ . We also assume that  $\Sigma_S$  is cycle free.

An example of  $\Sigma_S$  not being cycle-free would be when  $\Sigma_S(v_1) = v_2 \circ A_2$ , and  $\Sigma_S(v_2) = v_1 \circ A_1$ . Cycles inside  $\Sigma_C$  are unproblematic, see for example the definition of `fol` below in Example 1.9.

This completes the definition of the types, next we define adjectives:

**Definition 1.4.** We recursively define adjectives, starting with primitive adjective constructors:

- If  $c$  is a constant of one of the primitive types `bool`, `char`, `nat`, `double`, then  $c$  and  $c^\geq$  are adjectives.
- If  $c$  is a constant of type `selector`, then  $c$  is an adjective.
- **empty** is an adjective.
- If  $v$  is an identifier, then  $v$  is an adjective.

Next we define recursive adjective constructors:

- If  $f$  is an identifier, and  $A$  is an adjective, then  $f(A)$  is an adjective.
- If  $A$  is an adjective, then **first**( $A$ ) and **rest**( $A$ ) are adjectives.
- If  $A$  is an adjective, then  $\forall A$  and  $\exists A$  are adjectives.

Adjectives of any type can be combined by propositional operators:

- If  $A_1, \dots, A_n$  are adjectives, then  $A_1 \vee \dots \vee A_n$ , and  $A_1 \wedge \dots \wedge A_n$  are also adjectives.
- If  $A$  is an adjective,  $\neg A$  is also an adjective.

The intuitive meaning of the adjective  $c$  is : equal to  $c$ . The intuitive meaning of  $c^\geq$  is : greater or equal to  $c$ . The intuitive meaning of  $f(A)$  is: whose  $f$ -field satisfies  $A$ . If  $f$  is a repeated field, then  $A$  must be equal to **empty** or have one of the forms  $\forall A$ ,  $\exists A$ , **first**( $A$ ), or **rest**( $A$ ). The meaning of **empty** is that no repeated fields are present. The meaning of  $f(\forall A)$  is that all repetitions of  $f$  must satisfy  $A$ . Similarly  $f(\exists A)$  means that at least one repetition of  $f$  must satisfy  $A$ . The meaning of **first**( $A$ ) is the first repetition of  $f$  satisfies  $A$ . The meaning of **rest**( $A$ ) is that all repetitions of  $f$ , except for possibly the first, must satisfy  $A$ .

**Definition 1.5.** We assume a mapping  $\Sigma_A$  that maps identifiers to pairs of form  $(T, A)$ , where  $T$  is a simple type, and  $A$  is an adjective. If  $\Sigma_A$  contains a value for  $v$ , we write  $\Sigma_A(v)$  for the value.

The intuitive meaning of  $\Sigma_A(v) = (T, A)$  is: Identifier  $v$  is defined on type  $T$  as adjective  $A$ .

**Definition 1.6.** We write  $A[v]$  for an adjective that contains  $v$  somewhere not inside a subformula of form  $f(\cdot)$ . We define a cycle as a sequence of identifiers, s.t.

$$\Sigma_A(v_1) = (T_2, A_2[v_2]), \dots, \Sigma_A(v_{n-1}) = (T_n, A_n[v_n]), \Sigma_A(v_n) = (T_1, A_1[v_1]).$$

Cycles are logically problematic when they involve negation. Because the intended meaning of an adjective is an inductively defined predicate, an adjective defined through a cycle involving negation, would be ill-defined. This similar to the situation in logic programming ([7]). Although it would be possible to

adapt Definitions 2.2 and 3.7 to monotonic cycles, we are not aware of a meaningful use of it. Therefore, we think it is better to forbid cycles altogether. In the rest of this paper, we will assume that  $\Sigma_A$  is cycle free. We give some examples of type and adjective definitions:

**Example 1.7.** *The type of complex numbers can be defined as follows:*

$$\text{complex} := \text{re: double, im: double, } ( )^*$$

*Complex numbers have two scalar fields re and im, and no repeated fields.*

**Example 1.8.** *Natural numbers can be defined as follows:*

$$\text{nat} := \text{sel: selector, sel?} \left( \begin{array}{l} ?\text{zero} \Rightarrow ( )^* \\ ?\text{succ} \Rightarrow \text{pred: nat, } ( )^* \end{array} \right)$$

*Adjectives odd and even can be defined on nat in mutual recursion:*

$$\begin{aligned} \text{even} &:= ( \text{nat, sel}(?\text{zero}) \vee ( \text{sel}(?\text{succ}) \wedge \text{pred}(\text{odd}) ) ) \\ \text{odd} &:= ( \text{nat, sel}(?\text{succ}) \wedge \text{pred}(\text{even}) ) \end{aligned}$$

**Example 1.9.** *We define propositional logic:*

$$\text{prop} := \text{op: selector, op?} \left( \begin{array}{l} ?\text{var} \Rightarrow ( c: \text{char} )^* \\ ?\text{not} \Rightarrow \text{sub: prop, } ( )^* \\ ?\text{implies}\vee?\text{equiv} \Rightarrow \text{sub}_1: \text{prop, sub}_2: \text{prop, } ( )^* \\ ?\text{and}\vee?\text{or} \Rightarrow ( \text{sub}_n: \text{prop} )^* \end{array} \right)$$

There is no need to define  $\top$  or  $\perp$  separately, because  $\top = \bigwedge \emptyset$ , and  $\perp = \bigvee \emptyset$ . A variable consists of ?var combined with a finite number of characters. A negation consists of ?not combined with a single subformula called sub. A conjunction or disjunction consists of ?and or ?or, followed by an arbitrary number of subformulas called  $\text{sub}_n$ . We define a few adjectives on prop:

**Example 1.10.** *'Constructed by ?and' can be expressed as op(?and). 'Constructed by ?and or ?or' can be expressed as op(?and∨?or). The property of being an atom can be expressed by*

$$\text{atom} := ( \text{prop, op}(\text{var}) ).$$

*The property of being a literal can be expressed by*

$$\text{literal} := ( \text{prop, atom} \vee ( \text{op}(?\text{not}) \wedge \text{sub}(?\text{atom}) ) ).$$

*Negation normal form (NNF) can be expressed by:*

$$\text{nnf} := ( \text{prop, } \bigvee \left\{ \begin{array}{l} \text{literal} \\ \text{op}(?\text{and}\vee?\text{or}) \wedge \text{sub}_n(\forall \text{nnf}) \end{array} \right. ).$$

*Conjunctive normal form (CNF) can be defined as follows:*

$$\text{cnf} := ( \text{prop, op}(?\text{or}) \wedge \text{sub}_n(\forall(\text{op}(?\text{and}) \wedge \text{sub}_n(\forall \text{literal}))) ).$$

The sub field only exist when the op field equals ?not, which can be expressed by the adjective op(not). Similarly, field  $\text{sub}_1$  exists only when  $\text{op}(?\text{implies}\vee?\text{equiv})$ . The following definition defines the preconditions of field existence.

**Definition 1.11.** Let  $f$  be a exact identifier that is declared as field in a compound type  $C$ . We recursively define the precondition of  $f$ , written as  $\text{PREC}(C, f)$ , as follows:

- if  $C$  has form  $(v_1: V_1, \dots, v_n: V_n)^*$ , then  $\text{PREC}(f, C) = \top$ .
- If  $C$  has form  $v: V, C'$ , and  $f = v$ , then  $\text{PREC}(f, C) = \top$ . If  $f$  is declared in  $C'$ , then  $\text{PREC}(f, C) = \text{PREC}(f, C')$ .
- If  $C$  has form  $s?(A_1 \Rightarrow C_1, \dots, A_m \Rightarrow C_m)$ , and  $f = s$ , then  $\text{PREC}(f, C) = A_1 \vee \dots \vee A_m$ . Otherwise,  $f$  must be declared in exactly one  $C_j$ . We have  $\text{PREC}(f, C) = s(A_j) \wedge \text{PREC}(f, C_j)$ .

Since  $f$  can occur in  $\Sigma_C$  only once, it is possible to write  $\text{PREC}(f)$  instead of  $\text{PREC}(f, C)$ , where  $C$  is the compound type that contains  $f$ .

## 2. Semantics of Types and Adjectives

We will describe the semantics of adjectives in simplified fashion. We will assume that all data are represented by trees whose subtrees are labeled with field names. In reality, fields have a fixed position in memory, so that there is no need to store field names in data.

**Definition 2.1.** We define the set of data trees  $D$ , and the set  $\overline{D}$  of finite sequences of data trees, in simultaneous recursion:

- If  $d$  is an element of one of the primitive types  $T$ , defined in Definition 1.1, then  $d \in D$ .
- If  $d_1, \dots, d_n$  ( $n \geq 0$ ) is a finite sequence of data trees, then  $(d_1, \dots, d_n) \in \overline{D}$ . We will write  $\|\overline{D}\|$  for the length  $n$  of  $\overline{D}$ .
- If  $f_1, \dots, f_n$  are pairwise distinct identifiers, and each  $d_i \in D \cup \overline{D}$ , then  $\{(f_1, d_1), \dots, (f_n, d_n)\} \in D$ .

Data trees of the third type can be viewed as partial functions mapping identifiers to either a data tree or a sequence of data trees. We write  $d.f$  for the value attached to  $d$  assuming it exists. If  $d.f \in \overline{D}$ , we write  $d.f[i]$  for the  $i$ -th element of  $\overline{D}$ , assuming that  $i \leq \|d.f\|$ .

We will define when a data tree has a given type, and when a data tree (or sequence of data trees) satisfies a given adjective. Types may contain adjectives, but adjectives do not contain types. Therefore, we start by defining when a data tree (or sequence of data trees) makes an adjective true:

**Definition 2.2.** We define in simultaneous recursion when a data tree makes an adjective true, and when a sequence of data trees makes an adjective true. We use the notation  $d \vDash A$  for both cases.

We first consider the cases where  $d \in D$  :

- If  $c$  is a constant, then  $d \vDash c$  iff  $d = c$ , and  $d \vDash c^\geq$  iff  $d \geq c$ .
- If  $f$  is an identifier, s.t.  $d.f$  exists, then
  - if  $d.f$  is a single data tree, then  $d \vDash f(A)$  iff  $d.f \vDash A$ .
  - if  $d.f$  is a sequence of data trees,  $d \vDash f(A)$  iff  $d.f \vDash A$ .

The two cases appear to be the same, but  $d.f$  have different types, hence we prefer to list them separately.

- If  $v$  is an identifier, s.t.  $\Sigma_A(v)$  is defined, and  $d \in D$ , then  $d \vDash v$  iff  $d \vDash \Sigma_A(v)$ .

Next we list the cases where  $d \in \overline{D}$  :

- $(d_1, \dots, d_n) \vDash \text{empty}$  iff  $n = 0$ .
- $(d_1, \dots, d_n) \vDash \forall A$  iff for every  $d_i$  we have  $d_i \vDash A$ .
- $(d_1, \dots, d_n) \vDash \exists A$  iff there exists a  $d_i$  such that  $d_i \vDash A$ .
- $(d_1, \dots, d_n) \vDash \text{first}(A)$  iff  $n \geq 1$  and  $d_1 \vDash A$ .
- $(d_1, \dots, d_n) \vDash \text{rest}(A)$  iff  $n \geq 1$  and  $(d_2, \dots, d_n) \vDash A$ .

The definitions for the propositional connectives are standard, both for  $d \in D$  and  $d \in \overline{D}$  :

- $d \vDash A_1 \vee \dots \vee A_n$  if there is an  $i$  ( $1 \leq i \leq n$ ), s.t.  $d \vDash A_i$ .
- $d \vDash A_1 \wedge \dots \wedge A_n$  for all  $i$  ( $1 \leq i \leq n$ ), one has  $d \vDash A_i$ .
- $d \vDash \neg A$  iff not  $d \vDash A$ .

We define when a tree has a certain type:

**Definition 2.3.** We recursively define when a data tree  $d$  has type  $T$ . We use notation  $d: T$ . We first list the cases where  $T$  is simple:

- For a primitive type  $T$ , we define  $d: T$  as  $d \in T$ .
- If  $v$  is an identifier that is defined in  $\Sigma_S$ , then  $d: v$  iff  $d: \Sigma_S(v)$ .
- If  $v$  is an identifier that is defined in  $\Sigma_C$ , then  $d: v$  iff  $d: \Sigma_C(v)$ .
- $d: (T \circ A)$  iff  $d: T$  and  $d \vDash A$ .

Next we list the cases where  $T$  is compound:

- $d: (v_1: V_1, \dots, v_n: V_n)^*$  iff either
  - $n = 0$ , or
  - $n > 0$  and  $d.v_1, \dots, d.v_n$  are all defined, are all in  $\overline{D}$ , have the same length  $L = \|d.v_1\| = \dots = \|d.v_n\|$ , and for every  $i$  ( $1 \leq i \leq n$  and  $j$  ( $1 \leq j \leq L$ ), we have  $d.v_i[j]: V_i$ .
- $d: (v: V, C)$  iff  $d.v$  is defined,  $d.v: V$ , and  $d: C$ .
- $d: s?(A_1 \Rightarrow C_1, \dots, A_m \Rightarrow C_m)$  iff  $d.s$  is defined, in  $D$ , and there is exactly one  $1 \leq i \leq m$ , s.t.  $d.s \vDash A_i$  and  $d: C_i$ .

**Example 2.4.** Following up on Example 1.7, we can see that  $\{(\text{re}, 1.0), (\text{im}, 2.0)\}$  has type complex.

**Example 2.5.** Using the declaration of prop in Definition 1.9, the propositional formula  $p$  will be represented by a data tree  $d_p$  with fields

$$\begin{aligned} d_p.\text{op} &= \text{?atom} \\ d_p.\text{c} &= ('p') \end{aligned}$$

The propositional formula  $q$  will be represented by  $d_q$  with fields

$$\begin{aligned} d_q.\text{op} &= \text{?atom} \\ d_q.\text{c} &= ('q') \end{aligned}$$

Similarly, the propositional formula  $r$  will be represented by  $d_r$  with

$$\begin{aligned} d_r.\text{op} &= \text{?atom} \\ d_r.\text{c} &= ('r') \end{aligned}$$

The formula  $q \vee r$  will be represented by  $d$  with

$$\begin{aligned} d.\text{op} &= \text{?or} \\ d.\text{sub} &= (d_q, d_r) \end{aligned}$$

Finally, the formula  $p \rightarrow (q \vee r)$  will be represented by  $d'$  with

$$\begin{aligned} d.\text{op} &= \text{?implies} \\ d.\text{sub}_1 &= d_p \\ d.\text{sub}_2 &= d' \end{aligned}$$

### 3. A Tableaux Calculus

We will now define a calculus. As usual in theorem proving, the calculus checks satisfiability. This is sufficient to answer all logical questions that are needed for the implementation of our programming language. In order to show that adjectives  $A_1, \dots, A_n$  cover all possible cases in a switch, it is sufficient to show that  $A \wedge \neg A_1 \wedge \dots \wedge \neg A_n$  is unsatisfiable, where  $A$  is the adjective part of the type of the switch expression. In order to show that a switch has no overlapping cases, it is sufficient to show that  $A \wedge A_i \wedge A_j$  is unsatisfiable

for all distinct  $i$  and  $j$ . In order to show that a function  $f_1$  defined on adjective  $A_1$  is a better fit than function  $f_2$  defined on adjective  $A_2$ , it is sufficient to show that  $A_1 \wedge \neg A_2$  is unsatisfiable, while  $A_2 \wedge \neg A_1$  is satisfiable.

In order to use our calculus, types need to be decomposed into their adjective part, and the part that specifies their implementation. The calculus needs the adjective part of a type, but does not use the implementation. We will define two functions that decompose a simple type into their adjective and implementation components.

**Definition 3.1.** An implementation type is either a primitive type, or an identifier defined in  $\Sigma_C$ , i.e. the name of a compound type. For a simple type  $T$ , we define  $\text{ADJ}(T)$  and  $\text{IMPL}(T)$  as follows:

- If  $T$  is primitive then  $\text{ADJ}(T) = \top$ , and  $\text{IMPL}(T) = t$ .
- If  $v$  is an identifier defined in  $\Sigma_C$ , then  $\text{ADJ}(v) = \top$ , and  $\text{IMPL}(v) = v$ .
- $\text{ADJ}(T \circ A) = \text{ADJ}(T) \wedge A$ , and  $\text{IMPL}(T \circ A) = \text{IMPL}(T)$ .

For example, if  $T = \text{prop} \circ \text{nnf} \circ \text{cnf}$ , then  $\text{ADJ}(T) = \text{nnf} \wedge \text{cnf}$ , and  $\text{IMPL}(T) = \text{prop}$ .

**Theorem 3.2.** If  $T$  is a simple type, and  $d$  is a data tree, then  $d : T$  iff  $d : \text{IMPL}(T)$  and  $d \vDash \text{ADJ}(T)$ .

The theorem can be easily verified by applying the rules of Definition 2.2. At this point, we can simplify the interface of the tableaux procedure. Its input is just a single type  $T$ . The procedure establishes that there exists no data term  $d$  with type  $T$ . The applications that we mentioned in the introduction, can be obtained as follows:

- Establish that there exists no data term  $d$  of type  $T$  that satisfies adjective  $A$  : Call the procedure with  $T \circ A$ .
- Establish that every data term  $d$  of type  $T$ , that satisfies  $A$ , must also satisfy  $B$  : Call the procedure with  $T \circ A \circ \neg B$ .

The procedure works on the following normal form, which is similar to negation normal form:

**Definition 3.3.** An adjective is in path normal form (PNF) if negation is applied on identifiers only. An adjective  $A$  can be brought into PNF by calling  $\text{PNF}(A, \text{pos})$  defined below:

$$\begin{aligned}
\text{PNF}(A_1 \vee \dots \vee A_n, \text{pos}) &= \text{PNF}(A_1, \text{pos}) \vee \dots \vee \text{PNF}(A_n, \text{pos}) \\
\text{PNF}(A_1 \vee \dots \vee A_n, \text{neg}) &= \text{PNF}(A_1, \text{neg}) \wedge \dots \wedge \text{PNF}(A_n, \text{neg}) \\
\text{PNF}(A_1 \wedge \dots \wedge A_n, \text{pos}) &= \text{PNF}(A_1, \text{pos}) \wedge \dots \wedge \text{PNF}(A_n, \text{pos}) \\
\text{PNF}(A_1 \wedge \dots \wedge A_n, \text{neg}) &= \text{PNF}(A_1, \text{neg}) \vee \dots \vee \text{PNF}(A_n, \text{neg}) \\
\text{PNF}(f(A), p) &= f(\text{PNF}(A, p)) \text{ for } p \in \{\text{pos}, \text{neg}\} \\
\text{PNF}(\exists A, \text{pos}) &= \exists \text{PNF}(A, \text{pos}) \\
\text{PNF}(\exists A, \text{neg}) &= \forall \text{PNF}(A, \text{neg}) \\
\text{PNF}(\forall A, \text{pos}) &= \forall \text{PNF}(A, \text{pos}) \\
\text{PNF}(\forall A, \text{neg}) &= \exists \text{PNF}(A, \text{neg}) \\
\text{PNF}(\neg A, \text{pos}) &= \text{PNF}(A, \text{neg}) \\
\text{PNF}(\neg A, \text{neg}) &= \text{PNF}(A, \text{pos}) \\
\text{PNF}(v, \text{pos}) &= v \\
\text{PNF}(v, \text{neg}) &= \neg v
\end{aligned}$$

The following theorem can be straightforwardly proven, using Definition 2.2:

**Theorem 3.4.** Let  $A$  be an adjective. For every data tree  $d$ , we have  $d \vDash A$  iff  $d \vDash \text{PNF}(A, \text{pos})$ . Similarly,  $d \vDash \neg A$  iff  $d \vDash \text{PNF}(A, \text{neg})$ .

**Definition 3.5.** We define a path as a finite sequence  $(f_1, \dots, f_n)$  with  $(n \geq 0)$ , where each  $f_i$  is either an identifier (representing a field name), or an element of  $\{\exists, \forall, \text{first}, \text{rest}\}$ . We write  $\epsilon$  for the empty path, and use notation

$\pi.f$  for extending path  $\pi$  with  $f$ . If path  $\pi'$  can be obtained from path  $\pi$  by zero or more extensions, we call  $\pi$  a prefix of  $\pi'$ . We call  $\pi$  a strict prefix of  $\pi'$  if at least one extension was used.

We assume that there exists a total order  $<$  on paths with the property that if  $\pi_1$  is a strict prefix of  $\pi_2$ , then  $\pi_1 < \pi_2$ . Such an order can be easily obtained by fixing a total order on all possible  $f_i$ , and using the alphabetic, lexicographic extension.

**Definition 3.6.** An adjective stack  $S$  is a finite sequence of triples  $(\pi_i, A_i, \lambda_i)$ , where each  $\pi_i$  is a path, each  $A_i$  is an adjective, and each  $\lambda_i \in \mathcal{N}$ . We write  $S[\pi]$  for the set  $\{ A \mid \exists \lambda \text{ s.t. } (\pi, A, \lambda) \text{ occurs in } S \}$  and  $S_\lambda[\pi]$  for the set  $\{ A \mid (\pi, A, \lambda) \text{ occurs in } S \}$ .

The attribute  $\lambda_i$  stores the level of  $A_i$  on path  $\pi_i$ . More precisely, an adjective that was obtained from a formula on a strict prefix of  $\pi_i$  will have level 0. An adjective that was obtained from an adjective with level  $\lambda$  on the same path, receives level  $\lambda + 1$ .

We define the tableaux procedure. In order to obtain termination, the procedure uses a blocking rule. Whenever a new path  $\pi$  is entered, it checks that there is no strict prefix  $\pi'$  of  $\pi$  containing a set of adjectives that are included in the formulas on the current path. For example, if path  $(f)$  contains  $A, B$ , and path  $(f, g)$  contains formulas  $A, B, C$ , we can close the branch. This rule is correct, because for every data tree  $d$ , the subtree  $d.f.g$  is a subtree of  $d.f$ . If it is possible that  $d.f.g \models A \wedge B \wedge C$ , then one can replace  $d.f$  by  $d.f.g$  and obtain a smaller data tree, for which  $d.f \models A \wedge B$ .

The tableaux procedure always tries to extend in deterministic fashion first. If that does not result in a conflict, it selects a path  $\pi$  and tries all possible non-deterministic choices on this path. It keeps on trying until it either has explored all choices, or obtained a consistent stack  $S$ .

**Definition 3.7.** We define a procedure that tries to establish that no data term  $d$  can have type  $T$ . The procedure starts by creating a stack with one element:

$$S = ( (\epsilon, \text{ADJ}(T), 0) ).$$

After that, it calls  $b = \text{deterministic}(0)$ . If  $b = \perp$ , it returns  $\perp$ . Otherwise, it returns **nondeterministic**( $\epsilon$ ).

We list the subprocedures, starting with the procedure **deterministic**( $d$ ).

Procedure **deterministic**( $d$ ) must be called with a natural number  $1 \leq d \leq \|S\|$ . It returns  $\perp$  or  $\top$ , with  $\perp$  indicating contradiction. It is implemented as follows:

1. Set  $c = d$ .
2. If  $c \leq \|S\|$ , then
  - If  $A$  has form  $v$  with  $v$  an identifier that has no definition in  $\Sigma_A$ , then return  $\perp$ .
  - If  $A$  has form  $\neg c^\geq$  with  $c$  the minimal element in its type, then return  $\perp$ .
  - For all  $i < c$  with  $\pi_i = \pi_c$ , check whether  $A_i$  and  $A_c$  are in conflict, using the rules listed below. If they are in conflict, return  $\perp$ . Otherwise try the next  $\pi_i$ . The rules are:
    - A complementary pair  $A, \neg A$  is in conflict.
    - A pair  $c_1, c_2$  of any primitive type with  $c_1 \neq c_2$  is in conflict.
    - A pair of form  $\neg c_1^\geq, c_2^\geq$  with  $c_1 \leq c_2$  is in conflict.
    - A pair of form  $c_1^\geq, c_2$  with  $c_1 > c_2$ , is in conflict.
    - A pair of form  $\neg c_1^\geq, c_2$  with  $c_1 \leq c_2$  is in conflict.
  - Assign  $c = c + 1$  and go back to step 1.
3. If  $d \leq \|S\|$ , check if one of the deterministic extension rules in the table below is applicable on  $A_d$ . If yes, then for every deterministic consequence  $A'$ , push  $(\pi_d, A', \lambda_d + 1)$  to  $S$ .

$v$	$\Rightarrow$ PNF( $A$ , <b>pos</b> ) if $\Sigma_A$ contains a definition of form $v := (T, A)$
$\neg v$	$\Rightarrow$ PNF( $A$ , <b>neg</b> ) if $\Sigma_A$ contains a definition of form $v := (T, A)$
$\forall A$	$\Rightarrow$ <b>empty</b> $\vee$ ( <b>first</b> ( $A$ ) $\wedge$ <b>rest</b> ( $\forall A$ ))
$\exists A$	$\Rightarrow$ $\neg$ <b>empty</b> $\wedge$ ( <b>first</b> ( $A$ ) $\vee$ <b>rest</b> ( $\exists A$ ))
$A_1 \wedge \dots \wedge A_n$	$\Rightarrow$ $A_i$ for each $1 \leq i \leq n$



4. If  $A_d$  has form  $f(A')$ , then push  $(\pi_d.f, A', 0)$  and  $(\pi_d.f, \text{ADJ}(T), 0)$  to  $S$ , where  $\text{ADJ}(T)$  is the adjective component of the type  $T$  of  $f$ .

5. Set  $d = d + 1$ . Goto step 1.

Procedure **nondeterministic**( $\pi$ ) must be called with a path  $\pi$  that occurs in  $S$ . It returns  $\perp$  or  $\top$ , with  $\perp$  indicating contradiction. The implementation is as follows:

1. For every strict prefix  $\pi'$  of  $\pi$ , do:

(a) if  $S_0[\pi'] \subseteq S[\pi]$ , then return  $\perp$ . (This is the blocking rule, mentioned above.)

2. Call **nondeterministic**( $\pi, 1$ ).

Procedure **nondeterministic**( $\pi, n$ ) must be called with a path  $\pi$  that occurs in  $S$ , and with  $1 \leq n \leq \|S\|$ . It returns  $\perp$  or  $\top$ , with  $\perp$  indicating contradiction. The implementation is as follows:

1. Find the smallest  $n' \geq n$ , such that  $A_{n'}$  has form  $A_1 \vee \dots \vee A_m$  with  $m \geq 2$ .

2. If no  $n'$  was found in the previous step, then find the next path  $\pi' > \pi$  occurring in  $S$ , using the alphabetic lexicographic order  $>$ . Call  $b = \mathbf{nondeterministic}(\pi')$  and return  $b$ . If no path could be found, return  $\perp$ .

3. Set  $n = n'$ . We know that  $S_n$  has form  $(\pi, A_1 \vee \dots \vee A_m, \lambda)$ , with  $m \geq 2$ . For  $i$  from 1 to  $n$ , do the following:

- Set  $s = \|S\|$ . Push  $(\pi, A_i, \lambda + 1)$  on the stack.
- Call  $b = \mathbf{deterministic}(s)$ . If  $b \neq \perp$ , call  $b = \mathbf{nondeterministic}(\pi, n + 1)$ . If  $b = \top$  return  $\top$ .
- Restore  $S$  to length  $s$ .

4. Set  $n = n + 1$ . Restart at step 1.

Procedure **nondeterministic**( $\pi$ ) implements a loop checker which guarantees termination. Its correctness is based on the fact that, if a satisfying data tree can be found, it can be pruned into a data tree in which the calculus does not repeat initial states. We state the following without proof:

**Theorem 3.8.** *The tableaux calculus of Definition 3.7 terminates.*

In order to establish correctness of Theorem 3.8, it is sufficient establish that every branch is finite. During search, the procedure will only introduce subadjectives of the initial type, combined with subadjectives of identifiers defined in  $\Sigma_A$ . Therefore, blocking must eventually happen.

**Theorem 3.9.** *The procedure of Definition 3.7 is complete.*

*Доказательство.* The proof will be a bit informal. We have to prove that if no data tree  $d$  with  $d: T$  exists, then the procedure of Definition 3.7 will reject  $T$ . We will prove the converse: If  $T$  is not rejected, then there exists a data tree  $d$ , s.t.  $d: T$ . Suppose that  $T$  is not rejected. This means that there exists an open branch in the tableaux. Let  $S = (\pi_1, A_1, \lambda_1), \dots, (\pi_n, A_n, \lambda_n)$  be the stack with which the tableaux procedure terminated. For a path occurring in  $S$ , let

$$S_\pi = \{A \mid \text{there exists an } i \text{ s.t. } \pi_i = \pi \text{ and } A_i \text{ has form } c, c^\geq \text{ or } \neg c^\geq\}.$$

For a constant  $c$ , define

$$\begin{aligned} I(c) &= \{c\} \\ I(c^\geq) &= \{c' \mid c' \geq c\}, \\ I(\neg c^\geq) &= \{c' \mid c' < c\}. \end{aligned}$$

It is easily checked that for every  $\pi$  occurring in  $S$ , we have  $\bigcap \{I(A) \mid A \in S_\pi\} \neq \emptyset$ . Hence, we can select a constant  $c_\pi$  from each such set. Let  $A_S$  be the adjective

$$A_S = \bigwedge \{\pi(A_\pi) \mid \pi \text{ occurs in } S\}.$$

We prove by backwards induction (that is from  $n$  towards 1) It is easily checked that there exists a data tree  $d_S$  s.t.  $d_S \vDash A_S$ . We show by backwards induction that

$$d_S \vDash \pi_i(A_i) \wedge \dots \wedge \pi_n(A_n).$$

Assume that we already have  $d_S \vDash \pi_{i+1}(A_{i+1}) \wedge \dots \wedge \pi_n(A_n)$ . We proceed by case analysis on the form of  $A_i$ . We need to consider only the forms of  $A_i$  that do not close the branch.

- if  $A_i$  is an identifier  $v$ , then we know that  $v$  has a definition  $(T, A')$  in  $\Sigma_A$ , since otherwise the branch would have been closed. Hence we know that  $\text{PNF}(A', \text{pos})$  occurs among  $A_{i+1}, \dots, A_n$ . As a consequence,  $d_S \vDash \text{PNF}(A', \text{pos})$ . By Theorem 3.4, we know that  $d_S \vDash A'$ .
- If  $A_i$  is a negated identifier  $\neg v$ , then if  $v$  has no definition in  $\Sigma_A$ , we have  $d_S \not\vDash \pi_i(v)$ , so that  $d_S \vDash \pi_i(\neg v)$ . If  $v$  has a definition  $(T, A')$  in  $\Sigma_A$ , then  $\text{PNF}(A', \text{neg})$  occurs among  $A_{i+1}, \dots, A_n$ . By induction,  $d_S \vDash \text{PNF}(A', \text{neg})$ . By Theorem 3.4, we have  $d_S \vDash \neg A'$ .
- The other cases can be obtained by inspecting the cases in Definition 2.2.

□

**Theorem 3.10.** *The procedure of Definition 3.7 is sound.*

*Доказательство.* In order to prove soundness, one must prove that if the procedure rejects a type  $T$ , then there is no data tree  $d$  with  $d: T$ . We will prove the converse: If there is a data tree  $d$  with  $d: T$ , then  $T$  will not be rejected by the procedure.

This would be trivial, if the blocking rule would not exist. It is easy to show that if  $d: T$ , there exists a stack  $S = (\pi_1, A_1, \lambda_1), \dots, (\pi_n, A_n, \lambda_n)$ , with  $(\pi_1, A_1, \lambda_1) = (\epsilon, \text{ADJ}(T), 0)$ , that represents an open branch of the tableaux procedure when the blocking rule is not used.

Now assume that  $S$  will be closed when the blocking rule can be used. We will show that there exists a shorter stack  $S'$  which also starts with  $S'_1 = (\epsilon, \text{ADJ}(T), 0)$ , and which also represents an open branch of the tableaux procedure when the blocking rule cannot be used. Repeating this process will result in a stack that does not contain any more applications of the blocking rule.

Let  $S$  be a stack with  $S_1 = (\epsilon, \text{ADJ}(T), 0)$ , taken from an open branch of the tableaux procedure, when blocking is not used. Suppose that somewhere in  $S$ , the blocking rule would be applicable. This means that there exist  $\pi$  and  $\pi'$ , s.t.  $\pi'$  is a strict prefix of  $\pi$ , and  $S_0[\pi'] \subseteq S[\pi]$ . We prune  $S$  as follows:

- Remove every  $(\pi_i, A_i, \lambda_i)$ , s.t.  $\pi'$  is a prefix of  $\pi_i$  and  $\pi_i$  is a strict prefix of  $\pi$ .
- For every  $(\pi_i, A_i, \lambda_i)$ , s.t. which  $\pi$  is a prefix of  $\pi_i$  write  $\pi_i$  as  $\pi \cdot \pi'_i$ , and replace it by  $\pi' \cdot \pi'_i$ .

It can be checked that the resulting  $S'$  is shorter, because there is at least one  $(\pi_i, A_i, \lambda_i)$  with  $\pi_i = \pi'$  which will be removed. Moreover,  $S'$  still represents an open branch of the tableaux procedure without blocking. Hence we can continue the procedure and obtain an open branch on which the blocking rule is not applicable. □

#### 4. Resolving Overloads in Types and Adjectives

Until now we have insisted that field and adjective names are unique. In Example 1.9, we used  $\text{sub}$ ,  $\text{sub}_1$ ,  $\text{sub}_2$ , and  $\text{sub}_n$  as different variations of the name  $\text{sub}$ , dependent on whether we were taking a subformula of a negated formula, a formula constructed by a binary operator, or a formula constructed by an  $n$ -ary operator.

Similarly, we did not consider reuse of adjective names between different types. In reality, it is perfectly possible to have different types of formulas, for example modal, first-order and propositional, and define different nnf adjectives on each of them.

Modern programming languages like  $C^{++}$  or Java allow the use of ambiguous names which are made unambiguous by the compiler. For example, in  $C^{++}$ , one can define different print operators  $\ll$  on different types, and the compiler will pick the right one, when the programmer writes  $\ll$ . This is called *overload resolution*. Without overload resolution, the programmer needs to invent a new name for every type that needs to be printed. For example, in  $C$ , the programmer has to include the type in the name of a print function. This results in names like  $\text{printfol}$  or  $\text{printmodal}$ .

We want overload resolution in our programming language: It should be possible that the user defines different nnf adjectives on different types and reuses the same name 'nnf' for them. Similarly, we want that the user can call all variations  $\text{sub}$ ,  $\text{sub}_1$ ,  $\text{sub}_2$  and  $\text{sub}_n$ , just 'sub'.

Concretely, we want overload resolution on field names, adjective names, and function names. There will be no overload on type names, because we think it is unfeasible, and allowing ambiguous type names would result in ambiguous code.

In order to handle overload resolution, we introduce what we call *inexact identifiers*. Inexact identifiers are the identifiers that are used by the programmer in the program. We will call the identifiers that we have been using until now, *exact identifiers*. We assume a function  $v^?$  that maps exact identifiers  $v$  to their inexact representations. As an example, one can introduce an inexact name  $\text{sub}$  and set  $\text{sub}^? = \text{sub}_1^? = \text{sub}_2^? = \text{sub}_n^? = \text{sub}$ . Whenever identifier 'sub' occurs in the program, the compiler has to find out which of the exact variations of  $\text{sub}$  is meant.

**Definition 4.1.** *We assume an infinite set of inexact identifiers. We assume a map  $v^?$  that maps exact identifiers  $v$  to inexact identifiers.*

Now we explain how inexact identifiers are used in the program. In order to do that, we modify the definitions of type and of adjective. In Definition 1.2, in the definition of compound types, we make the following modifications:

- In case 1, the identifiers  $v_1, \dots, v_n$ , are inexact.
- In case 2, the identifier  $v$  is inexact.
- In case 3, identifier  $s$  is inexact.

In Definition 1.4, we will allow the identifiers  $v$  and  $f$  to be inexact. Apart from that, there are no changes.

**Example 4.2.** *We define propositional and multimodal logic, using inexact identifiers:*

$$\text{ident} := (\text{c} : \text{char})^*$$

$$\text{prop} := \text{op} : \text{selector}, \text{op}^? \left( \begin{array}{ll} ?\text{var} & \Rightarrow \text{ident}, ()^* \\ ?\text{not} & \Rightarrow \text{sub} : \text{prop}, ()^* \\ ?\text{implies}v? \text{equiv} & \Rightarrow \text{sub}_1 : \text{prop}, \text{sub}_2 : \text{prop}, ()^* \\ ?\text{and}v? \text{or} & \Rightarrow (\text{sub} : \text{prop})^* \end{array} \right)$$

Similarly, one can define multimodal logic with inexact identifiers:

$$\text{modal} := \text{op} : \text{selector}, \text{op}^? \left( \begin{array}{ll} ?\text{var} & \Rightarrow \text{ident}, ()^* \\ ?\text{not} & \Rightarrow \text{sub} : \text{prop}, ()^* \\ ?\text{implies}v? \text{equiv} & \Rightarrow \text{sub}_1 : \text{prop}, \text{sub}_2 : \text{prop}, ()^* \\ ?\text{and}v? \text{or} & \Rightarrow (\text{sub} : \text{prop})^* \\ ?\text{box}v? \text{dia} & \Rightarrow \text{sub} : \text{prop}, ()^* \end{array} \right)$$

In Example 4.2, the inexact field name  $\text{sub}$  occurs both in  $\text{prop}$  and in  $\text{modal}$ . In type  $\text{prop}$ , it occurs one time as scalar field, and one time as repeated field. If one knows that that  $f.\text{op} = ?\text{var}$ , one can access  $f.\text{sub}$ . If one knows that  $f.\text{op} \in \{?\text{and}, ?\text{or}\}$ , one can access  $f.\text{sub}[i]$ . Fieldname  $\text{sub}$  also occurs three times in type  $\text{modal}$ , two times as a scalar field, and one time as a repeated field.

Between the types  $\text{prop}$  and  $\text{modal}$ , we reuse the inexact identifiers  $\text{sub}_1$  and  $\text{sub}_2$ . It would not be possible to replace  $\text{sub}_1, \text{sub}_2$  be a single identifier, because they can occur at the same time in a formula formed by  $? \text{implies}$  or  $? \text{equiv}$ .

We define some ambiguous adjectives on  $\text{prop}$  and  $\text{modal}$  :

**Example 4.3.** *The adjectives atom and literal can be defined both on prop and on modal:*

$$\begin{aligned}
\text{atom: prop} & := \text{op}(\text{?var}) \\
\text{atom: modal} & := \text{op}(\text{?var}) \\
\text{literal: prop} & := \text{atom} \vee \text{op}(\text{?not}) \wedge \text{sub}(\text{?atom}) \\
\text{literal: modal} & := \text{atom} \vee \text{op}(\text{?not}) \wedge \text{sub}(\text{?atom})
\end{aligned}$$

*Similarly, NNF can be defined both on prop and on modal:*

$$\begin{aligned}
\text{nnf: prop} & := \bigvee \left\{ \begin{array}{l} \text{literal} \\ \text{op}(\text{?and}\vee\text{?or}) \wedge \text{sub}(\forall\text{nnf}) \end{array} \right. \\
\text{nnf: modal} & := \bigvee \left\{ \begin{array}{l} \text{literal} \\ \text{op}(\text{?and}\vee\text{?or}) \wedge \text{sub}(\forall\text{nnf}) \\ \text{op}(\text{?box}\vee\text{?dia}) \wedge \text{sub}(\text{nnf}) \end{array} \right.
\end{aligned}$$

In the example, there are different occurrences of sub, and it has to be determined which of the possible definitions is being referred to. Before we start discussing the treatment of inexact identifiers, note that defined identifiers in  $\Sigma_S$ ,  $\Sigma_C$  and  $\Sigma_A$  are always exact. Moreover, we do not allow the use of ambiguous type names. This means that when a defined or built-in type is used, it must be referred to by its exact name.

Resolving of inexact identifiers starts with a preprocessing step on  $\Sigma_S$  and  $\Sigma_C$ . During this step, declared fields are replaced by unique exact identifiers, and any adjectives used in field declarations are moved to  $\Sigma_A$ . This is done by replacing them with a unique exact identifier, and defining this identifier in  $\Sigma_A$ . After the preprocessing step, all inexact identifiers are in adjectives defined in  $\Sigma_A$ .

This replacement is somewhat similar to the way subformulas are replaced in the CNF transformation ([8]). In our case, the goal is not efficiency, but to make overload resolution possible in the first place. Adjectives occurring in the types of function declarations can be dealt with in the same way.

We allow reuse of the same inexact field name in a compound type, as long as the different occurrences of the field name are in different states of the type. For example in prop in Example 4.2, it is possible to reuse fieldname sub between formulas built by ?not and formulas built by ?and?or, but renaming both fields  $\text{sub}_1$  and  $\text{sub}_2$  into sub at the same time would be impossible. It would be still possible to rename one of them into sub.

**Definition 4.4.** *Let  $D = (w_1: W_1, \dots, w_n: W_n)$  be a sequence of declarations. Let  $w_{n+1}: W_{n+1}$  be a single declaration. We write*

$$D + (w_{n+1}: W_{n+1}) = (w_1: W_1, \dots, w_{n+1}: W_{n+1}).$$

**Definition 4.5.** *We define procedure  $\text{preprocsimple}(T)$  that preprocesses a simple type  $T$ . The result is again a simple type. It makes additions to  $\Sigma_A$  in the process.*

*The implementation of  $\text{preprocsimple}(T)$  is as follows:*

- Let  $A = \text{ADJ}(T)$ , and let  $T' = \text{IMPL}(T)$ . If  $A = \top$ , then return  $T'$ . If  $A \neq \top$ , then create a new exact identifier  $\alpha$ , assign  $\Sigma_A(\alpha) := (T', A)$ , and return  $T' \circ \alpha$ .

**Definition 4.6.** *We define procedure  $\text{preproccompound}(C, D)$  that preprocesses a compound type  $C$  within context  $D$ , where  $D$  is a sequence of declarations of form  $w: W$ , with  $w$  exact. The second argument  $D$  is used for checking that no inexact identifier occurs twice in the same path of  $C$ .*

- If  $C$  has form  $(v_1: V_1, \dots, v_n: V_n)^*$ , then
  - if there exist  $(w: W) \in D$  and  $1 \leq i \leq n$ , s.t.  $w^? = v_i$ , then create an error.
  - if there exist  $1 \leq i < j \leq n$ , s.t.  $v_i = v_j$ , then create an error.
  - Otherwise, create new, exact identifiers  $e_i$  for each  $v_i$ , and set  $e_i^? = v_i$ .

- Set  $V'_i = \mathbf{preprocsimple}(V_i)$ .
- Return  $(e_1: V'_1, \dots, e_n: V'_n)^*$ .
- If  $C$  has form  $(v: V, C')$  then if there is a  $(w: W) \in D$ , s.t.  $w^? = v$ , create an error.  
Otherwise, let  $e$  be a new, exact identifier for  $v$ . Set  $e^? = v$ . Set  $V' = \mathbf{preprocsimple}(V)$ . Return

$$e: V', \mathbf{preproccompound}(C', D + (e: V')).$$

- If  $C$  has form  $s?(A_1 \Rightarrow C_1, \dots, A_m \Rightarrow C_m)$ , then there must a be unique  $(w: W) \in D$ , s.t.  $w^? = s$ .
  - If no such  $w$  exists, or  $w$  is not unique, then create an error.
  - Otherwise, create new exact identifiers  $\alpha_1, \dots, \alpha_m$ , and set  $\Sigma_A(\alpha_j) = (W, A_j)$ , for each  $1 \leq j \leq m$ .
  - After that, return

$$w^?( \alpha_1 \Rightarrow \mathbf{preproccompound}(C_1, D), \dots, \alpha_m \Rightarrow \mathbf{preproccompound}(C_m, D) ).$$

We now apply procedures **preprocsimple** and **preproccompound** as follows:

- For every identifier  $v$  in the domain of  $\Sigma_S$ , replace  $\Sigma_S(v)$  by **preprocsimple**( $\Sigma_S(v)$ ).
- For every identifier  $v$  in the domain of  $\Sigma_C$ , replace  $\Sigma_C(v)$  by **preproccompound**( $\Sigma_C(v)$ , ( )).

**Example 4.7.** After being replaced by **preproccompound**, the definition of **prop** will have the following form in  $\Sigma_C$  :

$$\mathbf{prop} := \mathbf{op}: \mathbf{selector}, \begin{cases} \alpha_1 \Rightarrow \mathbf{ident}, ()^* \\ \alpha_2 \Rightarrow \mathbf{sub}: \mathbf{prop}, ()^* \\ \alpha_3 \Rightarrow \mathbf{sub}_1: \mathbf{prop}, \mathbf{sub}_2: \mathbf{prop}, ()^* \\ \alpha_4 \Rightarrow (\mathbf{sub}_n: \mathbf{prop})^* \end{cases}$$

$$\begin{aligned} \alpha_1: \mathbf{selector} &:= (\mathbf{selector}, ?\mathbf{var}) \\ \alpha_2: \mathbf{selector} &:= (\mathbf{selector}, ?\mathbf{not}) \\ \alpha_3: \mathbf{selector} &:= (\mathbf{selector}, ?\mathbf{implies}? \mathbf{equiv}) \\ \alpha_4: \mathbf{selector} &:= (\mathbf{selector}, ?\mathbf{and}? \mathbf{or}) \end{aligned}$$

At this point we have moved all inexactness into  $\Sigma_A$ . Adjective definitions in  $\Sigma_A$  have form  $\Sigma_A(v) = (T, A)$ , where  $v$  is an exact identifier,  $T$  is the type on which the defined adjective can be applied, and  $A$  is the adjective. Unfortunately,  $T$  is a simple type, which also may contain inexact adjectives. As an example where this could occur, one could define adjective **cnf** (Example 1.10) on **prop** $\circ$ **nfnf** instead of **prop**. The definition would have form  $\Sigma_A(\mathbf{cnf}) = (\mathbf{prop} \circ \mathbf{nfnf}, A')$ , where  $A'$  is an inexact version of the expression given in Example 1.10, and **nfnf** is an inexact identifier that must be resolved into the definition  $\Sigma_A(\mathbf{nfnf}) = (\mathbf{prop}, A'')$ , where  $A''$  is the expression of Example 4.3.

In order to solve this, we apply a preprocessing stage on  $\Sigma_A$ , similar to the preprocessing of  $\Sigma_S$ . We replace all adjectives occurring in domains by exact identifiers, while adding the definitions to  $\Sigma_A$  as follows:

- For every identifier  $v$  in the domain of  $\Sigma_A$ , write  $\Sigma_A(v)$  as  $(T, A)$ . Replace  $\Sigma_A(v)$  by **(preprocsimple**( $T$ ),  $A$ ).

Note that the call of **preprocsimple** may implicitly add new identifiers to  $\Sigma_A$ . This will not result in non-termination because these new identifiers need not be considered since their domain is an implementation type without adjectives. Unfortunately, there potentially exists a circular dependency between preconditions of adjectives that is hard to detect, as illustrated by the following example:

**Example 4.8.** Suppose that  $\Sigma_A$  contains adjective definitions with the following circular structure:

$$\begin{aligned} v_1 &:= (T_1 \circ \alpha_1, \dots) \\ v_2 &:= (T_2 \circ \alpha_2, \dots) \\ &\dots \\ \alpha_1 &:= (T_1, \dots, v_2^? \dots) \\ \alpha_2 &:= (T_2, \dots, v_1^? \dots) \end{aligned}$$

If at some other point an occurrence of  $v_1^?$  needs to be resolved, then  $v_1$  is an overload candidate. In order to decide if  $v_1$  should be used, one has to resolve the adjective of  $\Sigma_A(\alpha_1)$ , which contains  $v_2^?$ . In order to decide if  $v_2$  is an overload candidate for this occurrence of  $v_2^?$ , one has to resolve the adjective in  $\Sigma(\alpha_2)$ , which again contains  $v_1^?$ . This circular dependency is not solvable. Note that using  $v_1^?$  or  $v_2^?$  inside the right hand side of  $v_1$  or  $v_2$  is unproblematic.

The circular dependency in Example 4.8 is unsolvable, so the compiler has to reject it. Unfortunately, it is difficult to detect, because it only exists if the use of  $v_2^?$  in  $\Sigma(\alpha_1)$  is applied on implementation type  $T_2$ , and the use of  $v_1^?$  in  $\Sigma(\alpha_2)$  is applied on implementation type  $T_1$ . If for example the use of  $v_2^?$  in  $\Sigma(\alpha_1)$  is not on  $T_2$ , then  $v_2$  is not a candidate for overload, and there is no problem.

We solve this problem by starting to resolve an adjective definition, and whenever we encounter a precondition that has to be resolved first, we recursively try to solve this precondition. If this iteratively results in returning to the original definition, we reject  $\Sigma_A$ . In order to detect when we returned to the original definition, we maintain a set  $E$  of adjective definitions that we have encountered already. A circular dependency occurs when we need to resolve the precondition of an identifier that already occurs in  $E$ .

At this stage, all inexact identifiers are confined in  $\Sigma_A$ , in the second components  $A$  of the definitions  $\Sigma_A(v) = (T, A)$ . We give the procedure:

**Definition 4.9.** Let  $\Sigma_A$  be the map of inexact adjective definitions. We define procedure  $\text{RESOLVE}(v)$  that tries to resolve the overloads in  $\Sigma_A(v)$ .

It uses a set  $E$  of identifiers that were already encountered. Initially,  $E = \emptyset$ . The purpose of  $E$  is to detect circular dependencies of the type shown in Example 4.8. The implementation is as follows:

- If  $v$  has no definition in  $\Sigma_A(v)$ , then the result is an error.
- If  $v \in E$ , then the result is also an error. This means that a circular dependency was detected.
- Otherwise, add  $v$  to  $E$ .
- Write  $\Sigma_A(v)$  in the form  $(U, A)$ . If  $U$  has form  $(T \circ w)$ , then call  $\text{RESOLVE}(w)$ .
- Let  $A' = \text{RESOLVE}_T(\emptyset, A)$ , and replace  $\Sigma_A(v)$  by  $A'$ . Note that this function calls  $\text{RESOLVE}_T$ , that will be defined in Definition 4.10.

Function  $\text{RESOLVE}(v)$  possibly calls itself in order to make  $w$  exact. This will be detected, because we will have  $v \in E$ . There is no way to detect such circularity a priori, because its existence depends on the way overloads are resolved.

The second procedure  $\text{RESOLVE}_T(\Gamma, A)$  tries to resolve the fields and adjectives in the non-exact adjective  $A$  in context  $\Gamma$ . Context  $\Gamma$  is needed because it may contain preconditions of fields.

In case more than one overload candidate exists, we will take the nearest fit. This approach is used by  $C^{++}$  and Java. For example, if some class  $A$  is a subclass of  $B$ , which is a subclass of  $C$ , and some function  $f$  has definitions  $f(A)$  and  $f(C)$ , then the call  $f(a)$  will be resolved as  $f(A)$ , while the call  $f(b)$  will be resolved as  $f(C)$ . We have no notion of subclass, but we have implication between adjectives. Using implication between adjectives, the rule becomes as follows: If an application of some inexact identifier  $f$  has different overload candidates  $f_1, \dots, f_n$ , with  $f_1^? = \dots = f_n^? = f$ , s.t. each  $f_i$  is defined on adjective  $A_i$ , then we resolve  $f$  as follows: If there is a unique  $A_i$ , s.t.  $A_i$  implies all of  $A_1, \dots, A_n$ , then  $f$  will be resolved as  $f_i$ . In the definition, we will write  $\Gamma \vDash A$ , which means that  $\Gamma, \neg A$  is unsatisfiable. It should be noted that this is always in the context of a fixed  $\Sigma_S, \Sigma_C$ , and  $\Sigma_A$ . We do not want to use the notation  $\Sigma_S, \Sigma_C, \Sigma_A, \Gamma \vDash A$ , because it becomes too long.

**Definition 4.10.** Let  $T$  be an implementation type, let  $\Gamma$  be a set of exact adjectives applicable on  $T$ , and let  $A$  be an inexact adjective.  $\text{RESOLVE}_T(\Gamma, A)$  tries to resolve the overloads in  $A$  when applied on  $T$  in context  $\Gamma$ . If it succeeds it returns the exact version of  $A$ . We define  $\text{RESOLVE}_T(\Gamma, A)$  by cases on the form of  $A$ .

- If  $A$  is an (inexact) identifier  $v$ , then let  $v_1, \dots, v_n$  be the adjectives defined in  $\Sigma_A$ , that have  $v_i^? = v$  and for which  $\Sigma_A(v_i)$  has form  $(T, A_i)$  or  $(T \circ w_i, A_i)$ .  
Set  $C = \emptyset$ . For each  $i \in \{1, \dots, n\}$ , do the following:

- If  $w_i$  is absent, add  $i$  to  $C$ .
- Otherwise, call  $\text{RESOLVE}(w_i)$ . (This is the first version, defined in Definition 4.9.) If after that  $\Gamma \vDash w_i$ , then add  $i$  to  $C$ .

At this moment  $C$  is a subset of  $\{1, \dots, n\}$  containing the candidates that can still be considered as overloads for  $v$ . For  $i \in C$  do:

- for  $C \setminus \{i\}$  do:
  - \* If  $\Gamma, w_i \vDash w_j$ , then remove  $j$  from  $C$ .

If  $\|C\| \neq 1$ , then create an error message. Otherwise, return  $v_i$  where  $i$  is the unique element of  $C$ .

- If  $A$  is a constant  $c$ , then if  $c$  has primitive type  $T$ , return  $c$ . Otherwise create an error.
- If  $A$  has form  $c^\geq$ , and  $c$  has a primitive type that is not selector, then return  $c^\geq$ . Otherwise, create an error.
- If  $A$  has form  $f(A')$ , then let  $f_1, \dots, f_n$  be the fields (scalar or repeated) defined on type  $T$ , that have  $f_i^? = f$ . Note that  $T$  must be a defined compound type, because primitive types have no fields. Set  $C = \{1, \dots, n\}$ .

For each  $i \in C$  do:

- let  $\text{PREC}(f_i) = g_{i,1}(A_{i,1}) \wedge \dots \wedge g_{i,k_i}(A_{i,k_i})$  with  $k_i \geq 0$ , be the precondition of field  $f_i$  as defined in Definition 1.11. Each  $A_{i,j}$  is either a implementation type or has form  $T_{i,j} \circ w_{i,j}$  with  $w_{i,j}$  an exact identifier defined on implementation type  $T_{i,j}$ . For every  $w_{i,j}$  ( $1 \leq j \leq k_i$ ) that is present, call  $\text{RESOLVE}(w_{i,j})$ , defined in Definition 4.9.
- After that, check that  $\Gamma \vDash g_{i,1}(A_{i,1}) \wedge \dots \wedge g_{i,k_i}(A_{i,k_i})$ . If not, then remove  $i$  from  $C$ .

If  $\|C\| \neq 1$ , then the result is error. Otherwise, let  $i$  be the unique element of  $C$ . Assume that the declaration of  $f_i$  has form  $f_i: W$ . If  $f_i$  is a scalar field, then return

$$f_i(\text{RESOLVE}_{\text{IMPL}(W)}(\text{ADJ}(W), A')).$$

if  $f_i$  is a repeated field, then return

$$f_i(\text{RESOLVE}_{\text{IMPL}(W)}^*(\text{ADJ}(W), A')).$$

In the latter case, we called  $\text{RESOLVE}^*$  defined below.

- If  $A$  has form  $A_1 \vee \dots \vee A_n$ , then

$$\text{RESOLVE}_T(\Gamma, A) = \text{RESOLVE}_T(\Gamma, A_1) \vee \dots \vee \text{RESOLVE}_T(\Gamma, A_n).$$

- If  $A$  has form  $A_1 \wedge \dots \wedge A_n$ , then set  $\Gamma_1 = \Gamma$ . For  $i = 1$  to  $n$  do the following:
  - Let  $A'_i = \text{RESOLVE}_T(\Gamma_i, A_i)$ .
  - Set  $\Gamma_{i+1} = \Gamma_i \cup \{A'_i\}$ .

After that, return  $A'_1 \wedge \dots \wedge A'_n$ .

Next we define  $\text{RESOLVE}^*$  which resolves repeated fields.

- If  $A$  has form  $\forall A'$ , then return  $\forall \text{RESOLVE}_T(\emptyset, A')$ .
- If  $A$  has form  $\exists A'$ , then return  $\exists \text{RESOLVE}_T(\emptyset, A')$ .
- If  $A = \text{empty}$ , and  $T$  is not a compound type, then the result is an error. Otherwise return **empty**.

In order to make nnf of Example 4.3 exact, one has to start by calling  $\text{RESOLVE}(\text{nnf})$ . Procedure  $\text{RESOLVE}$  will insert nnf into  $E$ , and call  $\text{RESOLVE}_{\text{prop}}(\emptyset, A)$ , with  $A$  the expanded definition of nnf. Since  $A$  is a disjunction,  $\text{RESOLVE}_{\text{prop}}(\emptyset, A)$  will process the disjuncts independently.

The first disjunct equals literal. Procedure  $\text{RESOLVE}_{\text{prop}}(\emptyset, \text{literal})$  will establish that literal is the unique overload and return literal. It will not look at the definition of literal. If one wants to resolve the overloads in the definition of literal, one must call  $\text{RESOLVE}(\text{literal})$  separately.

The second disjunct equals  $\text{op}(\text{?andv?or}) \wedge \text{sub}(\forall \text{nnf})$ .  $\text{RESOLVE}_{\text{prop}}(\emptyset, \text{op}(\text{?andv?or}) \wedge \text{sub}(\forall \text{nnf}))$  will recursively call  $\text{RESOLVE}_{\text{prop}}(\text{op}(\text{?andv?or}))$ , which will resolve  $\text{op}$  (inexact) into  $\text{op}$  (exact) and call  $\text{RESOLVE}_{\text{selector}}(\text{andv?or})$ , which will return  $\text{?andv?or}$  since both are constants. After that, it will call

$$\text{RESOLVE}_{\text{prop}}(\text{op}(\text{?andv?or}), \text{sub}(\forall \text{nnf})).$$

There are two possible overloads for `sub`, namely `sub` for the `?not` case, and `subv`. We have  $\text{PREC}(\text{sub}) = \text{op}(\text{?not})$ , and  $\text{PREC}(\text{sub}_v) = \text{op}(\text{?andv?or})$ . Since only the latter is provable from the premiss, `subv` will be picked. Since `subv` is a repeated field, the procedure will recursively call  $\text{RESOLVE}_{\text{prop}}^*(\emptyset, \forall \text{nnf})$ , which will recursively call  $\text{RESOLVE}_{\text{prop}}(\emptyset, \text{nnf})$ . This call will return `nnf` without expanding it, after which the original call of  $\text{RESOLVE}_{\text{prop}}$  will construct the complete exact overload

$$\vee \left\{ \begin{array}{l} \text{literal} \\ \text{op}(\text{?andv?or}) \vee \text{sub}(\forall \text{nnf}) \end{array} \right.$$

## 5. Conclusions

Our goal is to implement an efficient programming language in which it is convenient to implement algorithms on trees whose forms are very different.

In order to obtain this, we have defined a flexible type system together with a way of refining these types by means of adjectives. The adjectives are intended as a replacement for matching in functional languages. In order to make this replacement possible, we have given a precise semantics for adjectives, so that adjectives can be evaluated on concrete data.

We provided a terminating tableaux calculus for deciding propositional relations between adjectives, and applied this procedure to overload resolution in imprecise formulations of adjectives. The overload resolution procedure replaces ambiguous field references in adjective definitions by exact field references. The same algorithm can be used for resolving ambiguous overloads in function calls.

## 6. Acknowledgements

This work gained from discussions with Cláudia Nalon. We thank Nazarbayev University for supporting this research through the Faculty Development Competitive Research Grant Program (FDCRGP) number 021220FD1651.

## Список литературы

- [1] Y. Minsky, A. Madhavapeddy и J. Hickey, *Real World OCaml (functional programming for the masses)*. O’Reilly, 2013, ISBN: 978-1-449-32391-2.
- [2] G. Hutton, *Programming in Haskell (second edition)*. Cambridge University Press, 2016, ISBN: 978-1-316-62622-1.
- [3] M. Odersky, L. Spoon и B. Venners, *Programming in Scala Third Edition*. Artima Press, 2007-2016, ISBN: 0-9815316-8-7.
- [4] P. Rondon, M. Kawaguchi и R. Jhala, “Liquid Types,” в *Programming Language Design and Implementation (PLDI)*, R. Gupta и S. Amarasinghe, ред., сер. ACM, ACM, 2009, с. 159–169.
- [5] R. Jhala и N. Vazou, “Refinement Types: A Tutorial,” *Foundations and Trends in Programming Languages*, т. 6, № 3–4, с. 159–317, 2021, ISSN: 2325-1107. DOI: [10.1561/2500000032](https://doi.org/10.1561/2500000032). url: <http://dx.doi.org/10.1561/2500000032>.
- [6] S. Klabnik и C. Nichols, *The Rust Programming Language*. No Starch Press, 2019, ISBN: 978-1718500440.
- [7] Michael Gelfond и Vladimar Lifschitz, “The Stable Model Semantics for Logic Programming,” в *Fifth International Conference and Symposium on Logic Programming*, R. Kowalski и K. Bowen, ред., MIT Press, 1988, с. 1070–1080.
- [8] G. Tseytin, “On the complexity of Derivation in the Propositional Calculus,” в *Studies in Constructive Mathematics and Mathematical Logic, Part II*, A. Slisenko, ред., Zapiski Nauchnykh Seminarov, 1970, с. 115–125.



---

## References

- [1] Y. Minsky, A. Madhavapeddy, and J. Hickey, *Real World OCaml (functional programming for the masses)*. O'Reilly, 2013, ISBN: 978-1-449-32391-2.
- [2] G. Hutton, *Programming in Haskell (second edition)*. Cambridge University Press, 2016, ISBN: 978-1-316-62622-1.
- [3] M. Odersky, L. Spoon, and B. Venners, *Programming in Scala Third Edition*. Artima Press, 2007-2016, ISBN: 0-9815316-8-7.
- [4] P. Rondon, M. Kawaguchi, and R. Jhala, "Liquid Types", in *Programming Language Design and Implementation (PLDI)*, R. Gupta and S. Amarasinghe, Eds., ser. ACM, ACM, 2009, pp. 159–169.
- [5] R. Jhala and N. Vazou, "Refinement Types: A Tutorial", *Foundations and Trends in Programming Languages*, vol. 6, no. 3–4, pp. 159–317, 2021, ISSN: 2325-1107. DOI: [10.1561/25000000032](https://doi.org/10.1561/25000000032). [Online]. Available: <http://dx.doi.org/10.1561/25000000032>.
- [6] S. Klabnik and C. Nichols, *The Rust Programming Language*. No Starch Press, 2019, ISBN: 978-1718500440.
- [7] Michael Gelfond and Vladimar Lifschitz, "The Stable Model Semantics for Logic Programming", in *Fifth International Conference and Symposium on Logic Programming*, R. Kowalski and K. Bowen, Eds., MIT Press, 1988, pp. 1070–1080.
- [8] G. Tseytin, "On the complexity of Derivation in the Propositional Calculus", in *Studies in Constructive Mathematics and Mathematical Logic, Part II*, A. Slisenko, Ed., Zapiski Nauchnykh Seminarov, 1970, pp. 115–125.