

# Implementing Set-Theoretic Types

MICKAËL LAURENT\*, Charles University, Czech Republic

KIM NGUYỄN\*, Université Paris-Saclay, CNRS, ENS Paris-Saclay, Laboratoire Méthodes Formelles, France

Set-theoretic types provide a rich type algebra that supports unrestricted unions, intersections, and negations, together with a decidable type constraint-solving algorithm known as tallying. These types are particularly well suited for typing dynamic languages, where functions often exhibit both generic and overloaded behavior. However, the complexity of their implementation has hindered their widespread adoption. In this paper, we introduce a modular representation for set-theoretic types and revisit the algorithms for subtyping and tallying. We compare our approach with the historical CDuce implementation and evaluate the performance impact of some optimizations and design choices.

CCS Concepts: • **Theory of computation** → *Type structures*; • **Software and its engineering** → *Polymorphism*; *Data types and structures*.

Additional Key Words and Phrases: set-theoretic types, semantic subtyping, tallying, implementation

## 1 Introduction

For over a decade, programmers have added static type systems to dynamic, originally untyped languages such as JavaScript and Python. These efforts aim to combine the flexibility of dynamic typing with the reliability of static type checking. Examples of such type systems are [TypeScript](#) and [Flow](#) for JavaScript or [Mypy](#), [Pyre](#) and [Pyright](#) for Python. To account for the variety of programming patterns expressible in untyped languages, these type systems often feature several (if not all) of the typing constructs that are part of the programming language literature, namely:

- parametric polymorphism (often referred to as *generics*)
- overloading or ad-hoc polymorphism, together with subtyping
- occurrence typing [28] (the refinement of a type following a dynamic type test)
- arbitrary union types (e.g. to express heterogeneous collections or optional results)
- intersection types (in particular for functions or records)
- gradual typing [27], the ability to mix typed and untyped code

Mixing together *all* of these, in a principled way, is a daunting task. Most, if not all the previously mentioned systems do not guarantee type safety and mix the various typing constructs in an ad-hoc way. Yet, for some of their use cases (documentation, test generation, code completion in IDE, ...) they provide excellent performance, and are able to infer *some* type information relatively quickly.

A theoretical formalism handling these features does exist: *set-theoretic types*. They were first introduced in the context of XML programming [18, 21] and later extended with parametric polymorphism [14], gradual typing [11] and type narrowing [12]. While set-theoretic types have started to be used as a theoretical foundation for type systems (e.g. [Etylizer](#) [26] for Erlang, [Elixir](#) [9]), they have not seen the wide adoption one could have hoped for. We believe that one of the main issues is the disconnect between the mathematically elegant theoretical foundation and the design and implementation issues a practical implementation faces.

The goal of this work is to report on SSTT, the simple set-theoretic type library, a reference implementation of set-theoretic types. We cover the implementation and data-structure of the basic type algebra, show how to extend it to cover practical (but often ignored) data types, and present

---

\*Both authors contributed equally to the paper

implementation techniques of high-level operators, namely subtyping and tallying (computing a set of substitutions that solve a subtyping constraint).

## 1.1 Overview

As a practical introduction to set-theoretic types, we use the following example from [14].

```

val map  :  $\forall \alpha, \beta. (\alpha \rightarrow \beta) \rightarrow [\alpha^*] \rightarrow [\beta^*]$ 
val even :  $\forall \gamma. (\text{int} \rightarrow \text{bool}) \wedge ((\gamma \setminus \text{int}) \rightarrow (\gamma \setminus \text{int}))$ 

```

Here, the `map` function has its familiar type, with the twist that list types are *regular expression types*. These are merely syntactic sugar for recursive types. For instance, the type  $[\alpha^*]$  is the recursive type defined by the equation  $X = \text{nil} \vee \alpha \times X$ , which is a *union* of `nil` – the *singleton* type of the constant representing the empty list – and a product type of the type variable  $\alpha$  (the polymorphic type of the head of the list) and  $X$  standing for the list type itself (the tail of the list). The `even` function is an *intersection* of two arrow types, meaning that `even` is an overloaded function. When this function is applied to a value of type `int`, it returns a value of type `bool`. Otherwise, when applied to values that are not integers, captured by the *set difference*  $\gamma \setminus \text{int}$ , it returns a value of that type<sup>1</sup>. Now, let’s imagine what would be needed to type the (partial) application `map even`. In a Hindley-Milner style type system, we would need to:

- (1) find a type substitution which unifies the type of the domain of `map` with the type of `even`
- (2) apply this substitution to  $[\alpha^*] \rightarrow [\beta^*]$  to deduce the type of the whole expression.

However, in the context of set-theoretic types, *syntactic* unification is not powerful enough. Indeed, it would fail here, since it confronts an arrow type  $(\alpha \rightarrow \beta)$  against an intersection type (the type of `even`). One of the key features of set-theoretic types is the associated notion of *semantic subtyping*. With this notion, the behavior of set-theoretic connectives w.r.t. type constructors (e.g. idempotence, distributivity) is derived from set theory. To reconcile parametric polymorphism and subtyping-based polymorphism, [14] introduces the *tallying* operation. In a nutshell, tallying consists in finding type substitutions for type variables that make a set of type inequations hold. Thus, typing the application `map even` consists in finding type substitutions such that

$$(\alpha \rightarrow \beta) \rightarrow [\alpha^*] \rightarrow [\beta^*] \leq ((\text{int} \rightarrow \text{bool}) \wedge ((\gamma \setminus \text{int}) \rightarrow (\gamma \setminus \text{int}))) \rightarrow \delta$$

where  $\leq$  denotes subtyping and  $\delta$  is a freshly introduced type variable that represents the type we are interested in, that is, the type of `map even`. This single problem contains the constraints that:

- the type of `even` must be a subtype of  $\alpha \rightarrow \beta$  (contravariance on the domain of `map`)
- $\delta$  must be a supertype of  $[\alpha^*] \rightarrow [\beta^*]$  (covariance on the codomain of `map`)

The result of this tallying problem is a set of four substitutions  $\{\sigma_1, \sigma_2, \sigma_3, \sigma_4\}$ , and the type of the application is therefore  $\delta\sigma_1 \wedge \delta\sigma_2 \wedge \delta\sigma_3 \wedge \delta\sigma_4$ , which becomes the type of the overloaded function:

$$\begin{aligned}
 & ([(\text{int} \vee \gamma_1)^*] \rightarrow [(\text{bool} \vee (\gamma_1 \setminus \text{int}))^*]) \wedge (\text{nil} \rightarrow \text{nil}) \wedge \\
 & ([(\gamma_3 \setminus \text{int})^*] \rightarrow [(\gamma_3 \setminus \text{int})^*]) \wedge ([\text{int}^*] \rightarrow [\text{bool}^*])
 \end{aligned}$$

This short example highlights most of the challenges one faces when trying to implement set-theoretic types. While some of them are documented in the literature, we propose in this work a systematic presentation of all the relevant aspects of the implementation of set-theoretic types, together with novel implementation techniques and insights discovered while implementing the SSTT library. Each technique is evaluated on a realistic use of set-theoretic types.

<sup>1</sup>To implement such behavior, the underlying programming language supports a dynamic type case construct.

*Contributions.* The contributions of this paper are as follows:

- We introduce *Binary Decision Trees* (BDTs), a data-structure used pervasively in our implementation of set-theoretic types. It is used to efficiently store and manipulate Boolean formulas involving custom atoms. We define a novel semantic simplification operation that reduces a BDT based on a semantic equivalence relation over the atoms (Section 3).
- We present a modular design for implementing set-theoretic types. We choose to represent them using a layered API, each layer handling a specific feature of set-theoretic types: type constructors, set-theoretic operators, type variables, recursivity (Section 4).
- We revisit the subtyping and tallying algorithms and propose new optimizations (Section 5).
- We present several useful extensions that can be encoded in the base type algebra (Section 6).
- We present our SSTT library, compare it against the CDuce implementation, and evaluate the impact of the semantic simplification and various optimizations (Section 7).

We recall the basics of set-theoretic types in Section 2 and review the related work in Section 8.

## 2 Set-theoretic types

### 2.1 Definitions

The types we implement in this paper are the set-theoretic types. For more detail, we refer the reader to the survey paper [8] which synthesises the relevant results from the literature.

*Definition 2.1 (Set-theoretic types).* The set  $\mathcal{T}$  of *set-theoretic types* is the set of regular and contractive terms coinductively defined by the following grammar:

$$\mathbf{Types} \quad t ::= b \mid \alpha \mid t \rightarrow t \mid t \times t \mid t \vee t \mid t \wedge t \mid \neg t \mid \mathbb{0} \mid \mathbb{1}$$

where  $b \in \mathcal{B}$  is a base type (in particular, it includes the constants  $c \in \mathcal{C}$  of the language we want to type) and  $\alpha \in \mathcal{V}$  is a type variable. The notation  $t_1 \setminus t_2$  is a syntactic sugar for  $t_1 \wedge \neg t_2$ . When writing a term, we use the following precedence (by decreasing priority):  $\neg, \setminus, \wedge, \vee, \times, \rightarrow$ .

As they are defined coinductively, types can be infinite trees, provided that they satisfy the constraints of regularity and contractivity explained below to ensure decidability of the subtyping relation. This yields a definition of equirecursive types that does not require explicit binders for recursion. A term is said *regular* if it only has a finite number of distinct subterms, and *contractive* if every infinite branch goes through an infinite number of arrows and products ( $\rightarrow$  and  $\times$ ).

The type  $\mathbb{0}$  is a special type that is not inhabited by any value, and is the subtype of all types. Conversely, the type  $\mathbb{1}$  is the supertype of all types. The  $\times$  constructor is used to type pairs, and the  $\rightarrow$  constructor is used to type functions. These are the most common type constructors, but set-theoretic types can be extended with other constructors as we will see throughout this paper.

Set-theoretic types are equipped with a decidable subtyping relation  $\leq$  (referred to as *semantic subtyping*, cf. Appendix A for more details). For this presentation, it suffices to consider that each type can be interpreted as a set of values that have that type, and that subtyping is set containment. Type connectives (union, intersection, negation) are interpreted as the corresponding set-theoretic operators. Although the interpretation of type variables is more complex, the following property is sufficient to grasp the behavior of subtyping in presence of type variables:

**PROPOSITION 2.2 (SUBTYPING, [16]).** *Let  $t$  and  $s$  be two types. The type  $s$  is a subtype of type  $t$  if, for every type substitution  $\sigma$ , we have  $\sigma s \leq \sigma t$ .*

We note  $\simeq$  the semantic equivalence:  $t_1 \simeq t_2$  if and only if  $t_1 \leq t_2$  and  $t_2 \leq t_1$ . Thanks to negation, the subtyping problem is equivalent to checking the emptiness of a type:  $s \leq t \iff s \setminus t \leq \mathbb{0}$ .

As presented in the introduction, subtyping is a building block for the definition of tallying:

*Definition 2.3 (Tallying, [14]).* Let  $S = \{(s_1, t_1), \dots, (s_n, t_n)\}$  be a finite set of pairs of types and  $\Delta$  a set of type variables. The solution of the tallying problem for  $S$  and  $\Delta$  is:

$$\text{tally}(S, \Delta) = \{\sigma \mid \text{dom}(\sigma) \cap \Delta = \emptyset \wedge \forall (s, t) \in S. s\sigma \leq t\sigma\}$$

The set  $S$  contains the subtyping constraints, and the set  $\Delta$  is the set of variables that the tallying is not allowed to instantiate. An important property of the tallying operation is that it is complete: it computes a finite set of substitutions which exactly characterize all the possible solutions.

## 2.2 Types in disjunctive normal forms

As hinted to in the previous section, deciding subtyping is equivalent to checking the emptiness of a type. A convenient way to express the subtyping algorithm is to put the type in *disjunctive normal form* (DNF). Given a type  $t$ , its DNF as defined in [16] is a syntactic term of the following form:

$$\text{DNF}(t) = \bigvee \left( \bigwedge_i \alpha_i^b \wedge \bigwedge_j \neg\beta_j^b \wedge \bigwedge_k b_k \wedge \bigwedge_l \neg b_l \right. \\ \bigvee \bigvee \bigwedge_i \alpha_i^p \wedge \bigwedge_j \neg\beta_j^p \wedge \bigwedge_k (t_k^1 \times t_k^2) \wedge \bigwedge_l \neg(s_l^1 \times s_l^2) \\ \left. \bigvee \bigvee \bigwedge_i \alpha_i^a \wedge \bigwedge_j \neg\beta_j^a \wedge \bigwedge_k (t_k^1 \rightarrow t_k^2) \wedge \bigwedge_l \neg(s_l^1 \rightarrow s_l^2) \right) \quad (*)$$

In this formula, each row is itself a DNF of positive and negative type variables and positive and negative instances of a particular type constructor: a basic type (we assume for now a single kind of basic type), a product, or an arrow. The subtyping algorithm can then iterate through the elements of this DNF to test the emptiness of the type. It is known that the complexity of semantic subtyping is in EXPTIME (see [20]). We attempt to tame this complexity in practice both at the level of type representation and in the subtyping and tallying algorithm.

## 3 Binary Decision Trees

In this section, we describe *Binary Decision Trees* (BDTs), a data structure for representing Boolean combinations of *atoms*. The nature of these atoms is varied, as shown in the previous section. They can be type variables, basic types, or type constructors. We develop the meta-theory of BDT and operations by abstracting over atoms. In addition to the standard Boolean operations on BDTs (Section 3.2), we define a semantic simplification operation (Section 3.3).

### 3.1 Structure of a BDT

Let  $\mathcal{A}$  be a set of atoms, ranged over by  $a$ . We assume we have a total order  $\preceq$  over atoms, as well as an interpretation  $\llbracket a \rrbracket$  of an atom  $a$  as a type. Let  $\mathcal{L}$  be a set of leaves, ranged over by  $l$ . We assume we have an interpretation  $\llbracket l \rrbracket$  of a leaf  $l$  as a type, a bottom leaf element  $\perp$  whose interpretation  $\llbracket \perp \rrbracket$  is the type  $\mathbb{0}$ , and a top leaf element  $\top$  whose interpretation  $\llbracket \top \rrbracket$  is the type  $\mathbb{1}$ . In addition, the set-theoretic operations  $\wedge$ ,  $\vee$ , and  $\neg$  must be defined on leaves, in accordance with the interpretation  $\llbracket \cdot \rrbracket$  ( $\llbracket l_1 \wedge l_2 \rrbracket \simeq \llbracket l_1 \rrbracket \wedge \llbracket l_2 \rrbracket$ ,  $\llbracket l_1 \vee l_2 \rrbracket \simeq \llbracket l_1 \rrbracket \vee \llbracket l_2 \rrbracket$ ,  $\llbracket \neg l \rrbracket \simeq \neg \llbracket l \rrbracket$ ).

*Definition 3.1 (BDT).* The set of Binary Decision Trees  $\text{BDT}(\mathcal{A}, \mathcal{L})$  over atoms  $\mathcal{A}$  and leaves  $\mathcal{L}$  are the finite trees produced by the following grammar (where  $l \in \mathcal{L}$  and  $a \in \mathcal{A}$ ):

$$\text{(Ordered) Binary Decision Trees } B ::= L(l) \mid N(a?B : B)$$

and satisfying the following properties:

(*ordering*) for any non-root node labeled  $a$  whose parent is labeled  $a'$ , we have  $a \not\preceq a'$ , and  
(*reduction*) all non leaf subtrees of the BDT are of the form  $N(a?B^+ : B^-)$  with  $B^+ \neq B^-$ .

In particular, the (*ordering*) property guarantees no branch goes through the same label twice ; the (*reduction*) property ensures a BDT does not represent a trivially empty type (cf. Definition 3.3).

Note that it is always possible to create BDTs that are reduced. To do so, we assume that  $N(\_? \_ : \_)$  behaves as a *smart constructor* which performs the following simplification:  $N(a?B : B) \rightsquigarrow B$ .

For concision, we use the notation  $\perp$  for  $L(\perp)$ ,  $\top$  for  $L(\top)$ , and  $N(a)$  for the BDT node  $N(a?\top : \perp)$ . The auxiliary definitions below allow us to extract the set of atoms or the set of leaves of a BDT:

*Definition 3.2.* For any BDT  $B$ , we define  $\text{leaves}(B)$  and  $\text{atoms}(B)$  as follows:

$$\begin{aligned} \text{leaves}(L(l)) &= \{l\} & \text{leaves}(N(a?B^+ : B^-)) &= \text{leaves}(B^+) \cup \text{leaves}(B^-) \\ \text{atoms}(L(l)) &= \emptyset & \text{atoms}(N(a?B^+ : B^-)) &= \text{atoms}(B^+) \cup \text{atoms}(B^-) \cup \{a\} \end{aligned}$$

The meaning of a BDT, that is, the type that it represents, is given by an interpretation function from BDTs to types.

*Definition 3.3.* The *interpretation*  $\llbracket B \rrbracket$  of a binary decision tree  $B$  is inductively defined as follows:

$$\llbracket L(l) \rrbracket = \llbracket l \rrbracket \quad \llbracket N(a?B^+ : B^-) \rrbracket = (\llbracket a \rrbracket \wedge \llbracket B^+ \rrbracket) \vee (\neg \llbracket a \rrbracket \wedge \llbracket B^- \rrbracket)$$

Note that the interpretation of a BDT  $B$  is a Boolean combination of atoms and leaves. An example of a BDT and the associated interpretation can be found in Figure 1. This Boolean combination can be expressed as a DNF (cf. Section 2.2):

*Definition 3.4.* The DNF  $\text{dnf}(B)$  of a BDT  $B$  is a set of triples  $(A_p, A_n, l)$ , where each triple represents a clause intersecting the positive atoms  $A_p$ , the negative atoms  $A_n$ , and the leaf  $l$ :

$$\begin{aligned} \text{dnf}(L(l)) &= \{(\emptyset, \emptyset, l)\} \\ \text{dnf}(N(a?B^+ : B^-)) &= \{(A_p \cup \{a\}, A_n, l) \mid (A_p, A_n, l) \in \text{dnf}(B^+)\} \cup \\ &\quad \{(A_p, A_n \cup \{a\}, l) \mid (A_p, A_n, l) \in \text{dnf}(B^-)\} \end{aligned}$$

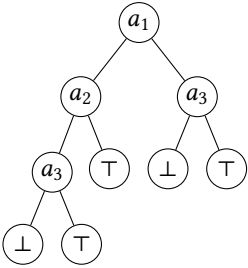


Fig. 1. BDT whose interpretation is  $\llbracket a_1 \rrbracket \wedge (\llbracket a_2 \rrbracket \wedge \neg \llbracket a_3 \rrbracket \vee \neg \llbracket a_2 \rrbracket) \vee \neg \llbracket a_1 \rrbracket \wedge \neg \llbracket a_3 \rrbracket$

The knowledgeable reader may recognize our BDTs as variants of Reduced Ordered Binary Decision Diagrams (ROBDD) [4]. There are however two key distinctions between the two formalisms. The first one is that for BDTs, leaf nodes may be arbitrary elements of a lattice (that can be mapped to the lattice of types). The second one is that we do not require BDTs to be DAGs, that is, we do not enforce that structurally equal subtrees are shared in memory. While it is possible to enforce this property (e.g. by performing hash-consing at construction time), it has two drawbacks. The first one is that, in typical use, it is seldom the case that

identical (and sufficiently large) subtrees are built to offset the cost of hash-consing. The second is that in our setting, we may have semantically equivalent subtrees which are not syntactically equivalent and thus will not be compacted by hash-consing. We set these considerations aside for now, and will revisit them in Section 7.

### 3.2 Set-theoretic operations

Implementing set-theoretic operations on BDTs is straightforward. For negation, it suffices to recursively traverse the BDT and negate the leaves:

*Definition 3.5.* The negation of a BDT can be computed by negating its leaves:

$$\neg L(l) = L(\neg l) \quad \neg N(a?B^+ : B^-) = N(a?\neg B^+ : \neg B^-)$$

For binary set-theoretic operations, their definitions can be given in a generic way by propagating the operations down the leaves, following the  $\preceq$  order over atoms.

*Definition 3.6.* Let  $\otimes \in \{\wedge, \vee, \setminus\}$  be a binary set-theoretic operation. Let  $B_1$  and  $B_2$  be two BDTs.  $B_1 \otimes B_2$  can be computed inductively as follows:

$$\begin{aligned} \mathsf{L}(l_1) \otimes \mathsf{L}(l_2) &= \mathsf{L}(l_1 \otimes l_2) \\ \mathsf{N}(a?B^+ : B^-) \otimes \mathsf{L}(l) &= \mathsf{N}(a?B^+ \otimes \mathsf{L}(l) : B^- \otimes \mathsf{L}(l)) \\ \mathsf{L}(l) \otimes \mathsf{N}(a?B^+ : B^-) &= \mathsf{N}(a?\mathsf{L}(l) \otimes B^+ : \mathsf{L}(l) \otimes B^-) \\ \mathsf{N}(a_1?B_1^+ : B_1^-) \otimes \mathsf{N}(a_2?B_2^+ : B_2^-) &= \begin{cases} \mathsf{N}(a_1?B_1^+ \otimes B_2^+ : B_1^- \otimes B_2^-) & \text{if } a_1 = a_2 \\ \mathsf{N}(a_1?B_1^+ \otimes \mathsf{N}(a_2?B_2^+ : B_2^-) : B_1^- \otimes \mathsf{N}(a_2?B_2^+ : B_2^-)) & \text{if } a_1 \preceq a_2 \\ \mathsf{N}(a_2?\mathsf{N}(a_1?B_1^+ : B_1^-) \otimes B_2^+ : \mathsf{N}(a_1?B_1^+ : B_1^-) \otimes B_2^-) & \text{if } a_2 \preceq a_1 \end{cases} \end{aligned}$$

Finally, we define a map operation on BDTs (it will be used later to implement type substitutions).

*Definition 3.7.* Let  $f : \mathcal{A} \cup \mathcal{L} \rightarrow \text{BDT}(\mathcal{A}, \mathcal{L})$  be a function from atoms and leaves to BDTs. We define the extension  $\bar{f} : \text{BDT}(\mathcal{A}, \mathcal{L}) \rightarrow \text{BDT}(\mathcal{A}, \mathcal{L})$  of  $f$  inductively, as follows:

$$\bar{f}(\mathsf{L}(l)) = f(l) \quad \bar{f}(\mathsf{N}(a?B^+ : B^-)) = (f(a) \wedge \bar{f}(B^+)) \vee (\neg f(a) \wedge \bar{f}(B^-))$$

### 3.3 Semantic simplification

Unlike BDDs traditionally used for representing Boolean formulas whose atoms are propositional variables, our BDTs may have more complex atoms (e.g. arrows, products). Two atoms may be semantically equivalent without being syntactically equivalent, or more generally, they may be correlated by the subtyping relation. These subtyping relations between the atoms is not captured by the total order  $\preceq$  we use for ordering the nodes of our BDTs, as this total order is purely syntactic. Hence, we define a simplification operation that removes nodes that are (semantically) redundant from a BDT. Note that this simplification operation does not try to reorder nodes or to change the atom they are labeled with. As such, it does not guarantee minimality of the resulting BDT, as a smaller semantically equivalent BDT that uses different atoms or a different order may exist.

We define our BDT semantic simplification operation inductively as follows:

$$\text{simpl}(t, \mathsf{L}(l)) = \mathsf{L}(l) \quad \text{simpl}(t, \mathsf{N}(a?B^+ : B^-)) = \begin{cases} B^{+'} & \text{if } t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \\ B^{-'} & \text{if } t \wedge \llbracket B^{-'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \\ B' & \text{otherwise} \end{cases}$$

where  $B^{+'} = \text{simpl}(t \wedge a, B^+)$ ,  $B^{-'} = \text{simpl}(t \wedge \neg a, B^-)$ , and  $B' = \mathsf{N}(a?B^{+'} : B^{-'})$ . The simplified form  $\text{simpl}(B)$  of a BDT  $B$  is then defined as  $\text{simpl}(B) = \text{simpl}(\mathbb{1}, B)$ .

In essence, the type  $t$  in  $\text{simpl}(t, B)$  represents the context of  $B$ , that is, the conjunction of positive and negative atoms traversed from the root to  $B$ . For each node  $B = \mathsf{N}(a?B^+ : B^-)$ , if we note  $t$  its context, the simplification procedure tests whether  $t \wedge \llbracket B \rrbracket \simeq t \wedge \llbracket B^+ \rrbracket$ . If that is the case, it means that the atom and right subtree are redundant and that  $B$  can be replaced by (the simplified version of)  $B^+$ . Otherwise, the same test is done for  $t \wedge \llbracket B \rrbracket \simeq t \wedge \llbracket B^- \rrbracket$ , which allows one to replace  $B$  by (the simplified version of) its right subtree. If both tests failed, the node is kept and its subtrees are recursively simplified.

Note that simplifying a BDT may be expensive as it requires deciding several semantic equivalences, each encoded as two subtyping tests. The simplification procedure is correct:

**PROPOSITION 3.8 (CORRECTNESS).** *For any BDT  $B$ , we have  $\llbracket \text{simpl}(B) \rrbracket \simeq \llbracket B \rrbracket$ .*

More interestingly, while the simplification procedure does not ensure minimality, it enjoys two desirable properties. First, the negation of a simplified type is itself simplified:

**PROPOSITION 3.9 (NEGATION).** *For any BDT  $B$ , we have  $\neg \text{simpl}(B) = \text{simpl}(\neg B)$  (where  $=$  is syntactic equality).*

Second, the simplification is sufficient to always reduce the empty type to the trivial BDT:

**PROPOSITION 3.10 (EMPTINESS).** *For any BDT  $B$ , we have  $\llbracket B \rrbracket \simeq \emptyset \Leftrightarrow \text{simpl}(B) = \perp$ .*

This guarantees that, if a BDT is simplified, then checking its semantic emptiness is equivalent to checking its syntactic emptiness. Both Propositions 3.9 and 3.10 will be exploited by our representation of set-theoretic types to improve performance (cf. Section 4.8). Proofs of these propositions are available in Appendix D.

## 4 Internal representation of types

In this section, we propose a structure to represent set-theoretic types using BDTs (cf. Section 3). This structure allows building types using different operations:

- type constructors and set-theoretic connectives ( $\wedge$ ,  $\vee$ ,  $\neg$ ),
- application of a type substitution  $\{\alpha_i \rightsquigarrow t_i\}_{i \in I}$  on a type  $t$ ,
- construction of recursive types from a set of equations  $\{\alpha_i = t_i\}_{i \in I}$ .

For simplicity, we only consider the constant, arrow, and product type constructors. Our structure is modular, so adding other type constructors is straightforward and is discussed in Section 6.

### 4.1 Overall structure

Set-theoretic types are defined coinductively and, thus, may be infinite trees. However, the regularity property ensures that the number of distinct subtrees is finite: we can thus represent a type as a graph that may contain cycles. We build this graph such that each subtree whose parent is a type constructor ( $\times$  or  $\rightarrow$ ) is represented by a node. A node thus defines the top-level structure of a subtree (i.e. the Boolean combination of constructors and type variables composing it), and references other nodes in the graph to describe the subtrees appearing inside a type constructor. The contractivity property of types ensures no such node definition consists in an infinite union or intersection. This graph structure is illustrated in Figure 2.

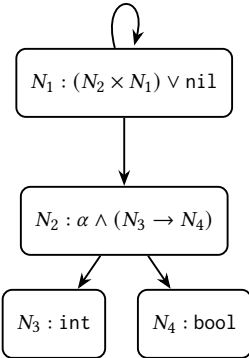


Fig. 2. Graph for the type  $\mu X. ((\alpha \wedge (\text{int} \rightarrow \text{bool})) \times X) \vee \text{nil}$

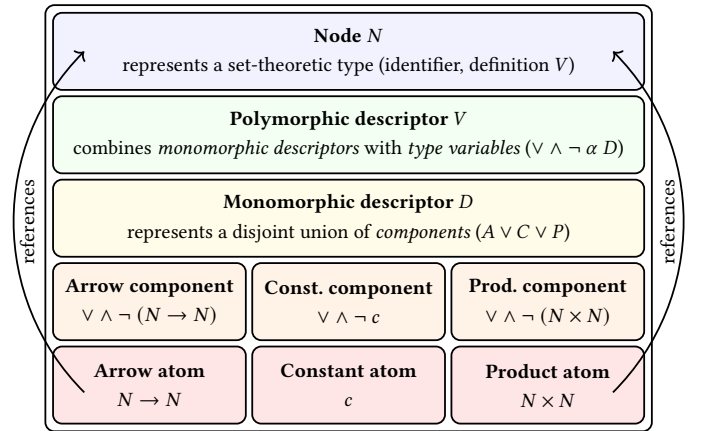


Fig. 3. Structure of a node (lower layers are used by upper layers)

Figure 3 illustrates the internal structure of each node of the graph. Each node can be decomposed into several layers. At the bottom, we have the different kinds of *atoms* (i.e. type constructors): products  $\times$ , arrows  $\rightarrow$ , and constants. Atoms of each kind are used inside *components* that represent a Boolean combination of atoms. Those different kinds of components are all disjoint (products, arrows and constants have disjoint interpretations): their disjoint union is represented by a *monomorphic descriptor*. Finally, a *polymorphic descriptor* expresses a Boolean combination of type variables and monomorphic descriptors.

In the following, nodes are ranged over by  $N$ . We assume we have a total order  $\preceq$  on nodes (it can be any total order, unrelated to the subtyping relation over types).

*Definition 4.1.* A *node substitution*  $\rho$  is a function from nodes to nodes which is the identity everywhere except for a finite set of nodes, called its domain and denoted by  $\text{dom}(\rho)$ .

## 4.2 Atoms

An atom  $a$  represents an instance of a type constructor. For an atom  $a$  of any kind (constant, product, arrow), we define its dependencies  $\text{deps}(a)$  and the application  $a\rho$  of a node substitution  $\rho$  to  $a$ .

*Arrow.* An arrow atom is a term  $N_1 \rightarrow N_2$ . Intuitively, it represents the type  $t_1 \rightarrow t_2$  where  $t_1$  (resp.  $t_2$ ) is the type associated with  $N_1$  (resp.  $N_2$ ).

$$\text{deps}(N_1 \rightarrow N_2) = \{N_1, N_2\} \quad (N_1 \rightarrow N_2)\rho = \rho(N_1) \rightarrow \rho(N_2)$$

*Product.* A product atom is a term  $N_1 \times N_2$ . Intuitively, it represents the type  $t_1 \times t_2$  where  $t_1$  (resp.  $t_2$ ) is the type associated with  $N_1$  (resp.  $N_2$ ).

$$\text{deps}(N_1 \times N_2) = \{N_1, N_2\} \quad (N_1 \times N_2)\rho = \rho(N_1) \times \rho(N_2)$$

*Constant.* A constant atom is a label  $c \in \mathcal{C}$  (e.g. true, false, nil).

$$\text{deps}(c) = \emptyset \quad c\rho = c$$

We also extend the total order  $\preceq$  over atoms (we can use the lexicographic order for arrows and products, and an arbitrary order over constants).

## 4.3 Components

A component represents a Boolean combination of atoms of a given kind. For simplicity, all kinds of components will be represented using BDTs. In practice, though, simple components can be represented using a more efficient, tailored representation (e.g. finite/cofinite sets for constants).

*Definition 4.2.* The BDTs we use for components use *Boolean leaves*, defined as follows:

$$\mathbf{Boolean\ Leaves} \quad b ::= \perp \mid \top$$

We define the interpretation of a Boolean leaf as follows:  $\llbracket \perp \rrbracket = 0$ ,  $\llbracket \top \rrbracket = 1$ . Set-theoretic operations ( $\wedge$ ,  $\vee$ ,  $\neg$ ) on Boolean leaves are defined accordingly.

*Definition 4.3.* A component  $C$  is a BDT whose leaves are Boolean leaves, and whose atoms are one kind of atoms defined above (arrows, products, constants).

Components inherit the operations defined on BDTs (cf. Section 3):  $\wedge$ ,  $\vee$ ,  $\neg$ ,  $\text{atoms}(C)$ ,  $\llbracket C \rrbracket$ ,  $\text{dnf}(C)$ . For each component  $C$  and node substitution  $\rho$ , we define  $\text{deps}(C)$  and  $C\rho$  as follows:

$$\text{deps}(C) = \bigcup_{a \in \text{atoms}(C)} \text{deps}(a) \quad C\rho = \tilde{f}(C)$$

where  $\tilde{f}$  is the extension on BDTs (cf. Definition 3.7) of the following function  $f$ :

$$f(b) = \mathbb{L}(b) \quad f(a) = \mathbb{N}(a\rho)$$

We also extend the total order  $\preceq$  over components (we can choose any syntactic total order over the structure of the underlying BDT).

#### 4.4 Monomorphic descriptors

A monomorphic descriptor  $D$  represents a disjoint union of components.

*Definition 4.4.* A *monomorphic descriptor* is a record  $\{\text{const}; \text{prod}; \text{arrow}\}$ , where  $\text{const}$  is a component for constants,  $\text{prod}$  is a component for products, and  $\text{arrow}$  is a component for arrows.

As these components have disjoint interpretations, all the set-theoretic operations ( $\neg$ ,  $\vee$ ,  $\wedge$ ) as well as node substitutions can be performed component-wise. The set of dependencies  $\text{deps}(D)$  of a descriptor  $D$  is the union of the dependencies of its components. We extend the total order  $\preceq$  over monomorphic descriptors (we can use the lexicographic order over the different components).

#### 4.5 Polymorphic descriptors

*Definition 4.5.* A *polymorphic descriptor*  $V$  is a BDT whose leaves are monomorphic descriptors, and whose atoms are type variables (we fix an arbitrary total order  $\preceq$  over type variables).

Polymorphic descriptors inherit the operations defined on BDTs (cf. Section 3):  $\wedge$ ,  $\vee$ ,  $\neg$ ,  $\text{atoms}(V)$ ,  $\text{leaves}(V)$ ,  $\llbracket V \rrbracket$ ,  $\text{dnf}(C)$ . We define  $\text{deps}(V)$  and  $V\rho$  (where  $\rho$  is a node substitution) as follows:

$$\text{deps}(V) = \bigcup_{D \in \text{leaves}(V)} \text{deps}(D) \qquad V\rho = \tilde{f}(V)$$

where  $\tilde{f}$  is the extension on BDTs (cf. Definition 3.7) of the following function  $f$ :

$$f(D) = \text{L}(D\rho) \qquad f(\alpha) = \text{N}(\alpha)$$

*Definition 4.6.* A *definition mapping*  $\sigma$  is a mapping from type variables to polymorphic descriptors. Its domain is denoted by  $\text{dom}(\sigma)$ .

The application of a mapping  $\sigma$  on a polymorphic descriptor  $V$  is defined as  $V\sigma = \tilde{f}(V)$ , where  $\tilde{f}$  is the extension on BDTs (cf. Definition 3.7) of the following function  $f$ :

$$f(D) = \text{L}(D) \qquad f(\alpha) = \begin{cases} \sigma(\alpha) & \text{if } \alpha \in \text{dom}(\sigma) \\ \text{N}(\alpha) & \text{otherwise} \end{cases}$$

We extend the total order  $\preceq$  over polymorphic descriptors (we can choose any syntactic total order over the structure of the underlying BDT).

#### 4.6 Nodes

*Definition 4.7.* A *node* is represented by a record  $\{\text{id}; \text{def}\}$  where  $\text{id}$  is a unique identifier (e.g. an integer), and  $\text{def}$  is a polymorphic descriptor.

The total order  $\preceq$  on nodes only depends on their  $\text{id}$  field. We write  $\text{node}(d)$  to denote a node with the definition  $d$  and a fresh identifier. Set-theoretic operations on nodes are defined as follows:

$$\neg N = \text{node}(\neg(N.\text{def})) \quad N_1 \wedge N_2 = \text{node}(N_1.\text{def} \wedge N_2.\text{def}) \quad N_1 \vee N_2 = \text{node}(N_1.\text{def} \vee N_2.\text{def})$$

The set of top-level type variables of a node  $N$  is defined as  $\text{top\_vars}(N) = \text{atoms}(N.\text{def})$ . The set of dependencies of a node  $N$ ,  $\text{deps}(N)$ , is the smallest set of nodes  $S$  such that  $N \in S$ , and  $\forall N' \in S. \text{deps}(N'.\text{def}) \subseteq S$ . The set of type variables of a node  $N$  is defined as  $\text{vars}(N) = \bigcup_{N' \in \text{deps}(N)} \text{top\_vars}(N')$ .

*Definition 4.8.* A *type substitution*  $\phi$  is a function from type variables to nodes which is the identity everywhere except for a finite set of type variables, denoted by  $\text{dom}(\phi)$ .

The application  $N\phi$  of a type substitution  $\phi$  on a node  $N$  is the node  $\rho(N)$  where:

- $\rho$  is a node substitution  $\{N_i \rightsquigarrow N'_i\}_{i \in I}$  where  $\{N_i\}_{i \in I} = \text{deps}(N)$ , and
- for every  $i \in I$ ,  $N'_i = \text{node}(((N_i.\text{def})\rho)\sigma)$ , and
- $\sigma$  is the definition mapping  $\{\alpha \rightsquigarrow \phi(\alpha).\text{def}\}_{\alpha \in \text{dom}(\phi)}$

Intuitively, the connected component of the node is copied ( $\rho$  associates each node to its copy), and the definition mapping  $\sigma$  is applied to the definition each copied node.

Note that this definition is recursive: the definition of the node  $N'_i$  depends on  $\rho$ , which itself involves  $N'_j$ . In SSTT, the field  $N.\text{def}$  is mutable, making it possible to generate a fresh node and only set its definition later (it also allows us to simplify the definition of a node after it is created, as it will be discussed in Section 4.8). Alternatively, these equations can be implemented by storing node definitions in a separate data-structure (for instance a dictionary mapping node identifiers to their definitions), or in a lazy language (e.g. Haskell) using recursive definitions.

#### 4.7 Construction of recursive types

Nodes can be constructed from atoms and combined using set-theoretic operations. However, we do not have a way to construct recursive types yet. We propose here a method for building a recursive type from a set of equations.

Let  $N$  be a node and  $E$  be a set of equations, each equation being a pair  $(\alpha, N_\alpha)$  corresponding to the equation  $\alpha = N_\alpha$ . Let  $S = \bigcup_{(\alpha, N_\alpha) \in E} \text{deps}(N_\alpha)$ . Using a topological sort algorithm, we order  $S$  into an ordered set  $\{N_i\}_{i \in 1..n}$  such that  $\forall i \in 1..n. \forall \alpha \in \text{top\_vars}(N_i). \forall j \in i..n. (\alpha, N_j) \notin E$ .

In other words, our order must guarantee that for every binding  $(\alpha, N_\alpha)$  of our set of equations  $E$ ,  $N_\alpha$  is smaller than all the nodes that depend on  $\alpha$  at top-level. If no such ordering exist, then it means that the set of equations is not contractive (it contains dependency cycles that do not pass under a type constructor): there may be no solution or infinitely many solutions that satisfy the set of equations, and thus we reject it.

Otherwise, for every  $i \in 1..n$ , we define a node  $N'_i = \text{node}(((N_i.\text{def})\rho)\sigma_i)$ , where:

- $\rho$  is the node substitution  $\{N_j \rightsquigarrow N'_j\}_{j \in 1..n}$ , and
- $\sigma_i$  is the definition mapping  $\{\alpha \rightsquigarrow N'_j.\text{def} \mid j \in 1..(i-1), \alpha \in \mathcal{V} \text{ s.t. } (\alpha, N_j) \in E\}$

As for type substitutions, this definition is recursive and should be implemented using similar techniques. The function  $\text{build}(N, E)$  returns a node  $N\phi$  where  $\phi = \{\alpha \rightsquigarrow N'_i \mid (\alpha, N_i) \in E\}$ . If  $E$  is not contractive,  $\text{build}(N, E)$  is not defined.

#### 4.8 Node simplification

Depending on how they are constructed, our nodes may contain redundant atoms in their internal representation. For instance, consider the following sequence of operations: (i) the type  $\text{bool} \rightarrow \text{bool}$  is constructed, (ii) it is intersected with  $\text{true} \rightarrow \text{false}$  (where  $\text{true}$  and  $\text{false}$  are singleton types that are subtypes of  $\text{bool}$ ), and (iii) it is intersected with  $\text{false} \rightarrow \text{true}$ . The type we obtain is semantically equivalent to  $(\text{true} \rightarrow \text{false}) \wedge (\text{false} \rightarrow \text{true})$ , but its internal representation will still contain the atom  $\text{bool} \rightarrow \text{bool}$ . Redundancy can also accumulate after applying multiple successive substitutions on a type, which typically happens when implementing a type inference algorithm based on tallying such as [12].

The BDT simplification procedure  $\text{simpl}(\cdot)$  defined in Section 3.3 can be used to avoid the accumulation of redundant atoms. In SSTT, we implemented a systematic simplification of types: a node is automatically simplified after being constructed by a  $\wedge$ ,  $\vee$ , a substitution, or a build operation. According to Proposition 3.9, a type does not need to be simplified again after a  $\neg$  operation. Keeping node definitions in a simplified form allows us to decide emptiness just by checking if the definition of the node is syntactically the empty BDT (cf. Proposition 3.10). This

means that, though the subtyping algorithm (Section 5) is used to simplify the definition of a node initially, subsequent emptiness checks are constant time. The time overhead induced by a systematic simplification of types is measured and discussed in Section 7.

## 5 Type operations

In the rest of this paper, we assimilate a type  $t$  with a node  $N$  that represents it. Set-theoretic operations and substitutions on types correspond to the associated operations on nodes.

### 5.1 Subtyping

As already explained, the subtyping algorithm is really an emptiness test. Following [16], to test a type  $t$  for emptiness one only has to:

- Compute  $\text{DNF}(t)$  as in equation (\*)
- Disregard the top-level variables  $\alpha_i^{\{\text{b,p,a}\}}$  and  $\beta_j^{\{\text{b,p,a}\}}$
- Test that all combinations of basic types, product types and arrow types are empty

The reason for the second condition is that, since the type must be empty *for all* possible substitutions, in particular for those mapping positive variables to  $\mathbb{1}$  and negative variables to  $\mathbb{0}$ . We detail a first version of the emptiness test in Algorithm 1.

ALGORITHM 1 (SUBTYPING, NON-OPTIMIZED).

```

1 let rec is_empty N Σ =
2   if N ∈ Σ then true
3   else
4     let Σ' = Σ ∪ {N} in
5     let V = N.def in
6     ∀D ∈ leaves(V),
7       is_empty_const D.const Σ' ∧
8       is_empty_prod D.prod Σ' ∧
9       is_empty_arrow D.arrow Σ'

```

In this code<sup>2</sup>,  $N$  is the node of the type we are testing and  $\Sigma$  is a set of nodes. Each layer (in the sense of Figure 3) handles a particular aspect of the emptiness test. Recursive types are handled at the level of nodes, by remembering which node has already been visited. If a node is encountered during its own traversal, we return that it is empty, since a recursive type whose emptiness depends only on itself, e.g.  $\mu X.(\mathbb{1}, X)$ <sup>3</sup>, is empty. The first time a node is visited, it is recorded (l. 4), and its polymorphic descriptor  $V$  (a BDT whose atoms are type variables) is inspected. At that point, it is sufficient to inspect all the monomorphic descriptors at the leaves of the BDT. Each such descriptor  $D$  is a disjoint union of BDTs of constants, products, and arrows, for which the emptiness test may call `is_empty` again on nested nodes. Thus, the two base cases of this recursive algorithm are: when a node has already been visited (l. 2), which will happen if the input type is recursive, and when the definition of the node is trivial and none of the calls to `is_empty_const`, `is_empty_prod`, and `is_empty_arrow` (l. 7-9) recurses. These functions are detailed below.

*Constant component.* Since constants are disjoint one from the other, testing the emptiness of a BDT of constants is done in constant time, by checking if it is structurally equal to the leaf  $\perp$ .

*Product component.* Testing the emptiness of products is done by the `is_empty_prod` function:

```
let is_empty_prod p Σ = ∀ ({P1i × P2i}i∈I, An, T) ∈ dnf(p). Ψprod(∧i∈I P1i × ∧i∈I P2i, An, Σ)
```

This function enumerates each conjunction of the DNF of  $p$  whose leaf node is  $\top$ . For each one, the positive part of the intersection is transformed into a single product, by pushing the intersection

<sup>2</sup>We use an OCaml-like syntax to present algorithm.

<sup>3</sup>The model of semantic subtyping only allows for finite values.

below the product constructor. It then remains to test whether the negative part of the intersection is enough to negate this single product, which is tested by the  $\Psi_{\text{prod}}$  function:

$$\begin{aligned}\Psi_{\text{prod}}(P_1 \times P_2, \emptyset, \Sigma) &= (\text{is\_empty } P_1 \Sigma) \vee (\text{is\_empty } P_2 \Sigma) \\ \Psi_{\text{prod}}(P_1 \times P_2, \{N_1 \times N_2\} \cup S, \Sigma) &= (\text{is\_empty } P_1 \Sigma) \vee (\text{is\_empty } P_2 \Sigma) \vee \\ &\quad (\Psi_{\text{prod}}((P_1 \setminus N_1) \times P_2, S, \Sigma) \wedge \Psi_{\text{prod}}(P_1 \times (P_2 \setminus N_2), S, \Sigma))\end{aligned}$$

Note that in the worst case, the  $\Psi_{\text{prod}}$  function may perform an exponential number of calls to `is_empty`, each Boolean connective introducing a potential backtracking point. Crucially, even though new types are generated (e.g. when pushing the intersection below the products or when computing type differences), only a finite number of such new types are created (see [16, Section 3.5]) which ensures the termination of the algorithm (the number of all nodes collected in  $\Sigma$  is finite).

*Arrow component.* Testing the emptiness of arrows works similarly to the case of products. The only difference is that intersection and arrow constructors do not commute, and that a positive intersection of arrows is never empty. We can apply these principles to test emptiness as follows.

**let** `is_empty_arrow`  $a \Sigma = \forall (A_p, A_n, \top) \in \text{dnf}(a). \exists (N_1 \rightarrow N_2) \in A_n. \Psi_{\text{arrow}}(N_1, \neg N_2, A_p, \Sigma)$

For an intersection of positive arrows  $A_p$  and negative arrows  $A_n$ , there must be a single negative arrow  $N_1 \rightarrow N_2$ , which negates the whole positive intersection. For that to be true, it suffices that at least one arrow in  $A_p$  is completely negated by the selected negative arrow (which makes the positive intersection empty). This is done by the auxiliary  $\Psi_{\text{arrow}}$  function defined as:

$$\begin{aligned}\Psi_{\text{arrow}}(N_1, T_2, \emptyset, \Sigma) &= (\text{is\_empty } N_1 \Sigma) \vee (\text{is\_empty } T_2 \Sigma) \\ \Psi_{\text{arrow}}(N_1, T_2, \{P_1 \rightarrow P_2\} \cup A_p, \Sigma) &= (\text{is\_empty } N_1 \Sigma) \vee (\text{is\_empty } T_2 \Sigma) \vee \\ &\quad (\Psi_{\text{arrow}}((N_1 \setminus P_1), T_2, A_p, \Sigma) \wedge \Psi_{\text{arrow}}(N_1, (P_2 \wedge T_2), A_p, \Sigma))\end{aligned}$$

This function is similar to the one used for products, but takes a single negative arrow  $N_1 \rightarrow N_2$  and check that it is a super type of the positive part (and thus, that their difference is empty). Notice that  $T_2$  is initially the negation of the co-domain  $N_2$ .

*Recursive types and caching.* To explain how recursive types are handled, consider the example:

$$X = (Y \times Y) \quad Y = (\text{nil} \times Y) \vee (Z \times Z) \vee (\text{nil} \rightarrow \text{nil}) \quad Z = Y \times \text{nil}$$

where types are written as a set of recursive equations between nodes (and some basic types). If one uses the algorithm of Algorithm 1 to test the emptiness of  $X$ , the result is `false`. A relevant subset of the recursive calls is given in Figure 4. To test the emptiness of  $X$ , we must test its definition,

```

1 is_empty X ∅
2   is_empty_prod Y × Y {X}
3     is_empty Y {X}
4       is_empty_prod (nil × Y) ∨ (Z × Z) {X, Y}
5         (is_empty_const nil {X, Y} ∼ false
6           ∨ is_empty Y {X, Y} ∼ true) (*recursive occurrence*)
7         ∧ is_empty_pair Z × Z {X, Y, Z} ∼ true (*depends on Y*)
8       ∧ is_empty_arrow (nil → nil) {X, Y} ∼ false
9     ∨ is_empty Y {X}

```

Fig. 4. Recursive calls for checking testing the emptiness of  $X$

since  $X$  is not in  $\Sigma$  (l. 2). Its definition is the product type  $Y \times Y$ , which is empty if either projection is empty. Those are tested at l. 3 and l. 9. This illustrates an inefficiency of the algorithm. The set  $\Sigma$  is used to handle recursion, and not as a cache, which means that several occurrences of the same node may be traversed several times.

While it is tempting to replace sigma with a mutable map, one must take special care in doing so. Indeed, in this example, it would be recorded that:

- at first,  $Y$  is empty (initial recursive call, l. 3)
- during its traversal,  $Z$  (which depends on  $Y$ ) is found to be empty
- ultimately,  $Y$  contains  $\text{nil} \rightarrow \text{nil}$  and therefore is *not* empty

It is not sufficient to update  $\Sigma$  to record that  $Y$  is non-empty. One must also *invalidate* the results stored for all nodes which used the *wrong* initial assumption that  $Y$  was empty. This improvement is described in Algorithm 2. The algorithm maintains three data-structures:

ALGORITHM 2 (SUBTYPING, OPTIMIZED).

```

1 let memo = H.create ()
2 let stack = ref []
3 let mem_stack = H.create ()
4 let rec is_empty N =
5   if N ∈ memo then begin
6     if N ∈ mem_stack then
7       H.add mem_stack N
8       (!stack::H.find mem_stack N);
9       H.find memo N
10  end else begin
11    stack := N :: !stack;
12    H.add mem_stack N [];
13    H.add memo N true;
14    let V = N.def in
15    let b = ∀D ∈ leaves(V),
16      is_empty_const D.const ∧
17      is_empty_prod D.prod ∧
18      is_empty_arrow D.arrow in
19    if not b then
20      invalidate N memo mem_stack;
21    H.add memo N b;
22    H.remove mem_stack N;
23    stack := List.tl !stack;
24    b end

```

- the hash table memo (l. 1)
- a reference to a persistent list, used as a stack (stack, l. 2)
- a hash table mapping nodes to stacks (mem\_stack, l. 3)

While memo plays the same purpose as in the previous algorithm, the other two data-structures are used to track dependencies between nodes. Assuming a node  $N$  is not in memo:

- push it on stack (l. 11)
- create an entry in mem\_stack (l. 12)
- map it to true in memo as usual (l. 13)

We then explore the descriptor of the node recursively. If a node  $N$  is in memo, then:

- if an entry exists for  $N$  in mem\_stack, record the current stack (l. 6–8)
- return the memoized value for  $N$  (l. 9)

Lastly, after returning from a recursive call, if the node turns out to be non-empty then:

- for each recorded stack for  $N$ , remove from memo any node  $X$  appearing on that stack, until  $N$  is reached (done by function invalidate, l. 20)
- update memo with the result for  $N$  (l. 21)

- remove  $N$  from mem\_stack (l. 22)
- pop  $N$  from the stack (l. 23)
- return the result (l. 24)

Whenever we encounter a node  $N$  a second time during a recursive traversal, either:

- it is in a definitive state (it has no binding in mem\_stack), we can therefore directly return its emptiness state from the memo table;
- it is still in an “unknown” state (i.e. it has a binding in mem\_stack), and every node that has been put on the stack during its inspection depends on the (possibly wrong) assumption that it is empty, therefore we must remember them to remove them from memo later on.

## 5.2 Tallying

We revisit the tallying algorithm formalized by [14], which does not perform well in practice. At its heart, the tallying algorithm is fairly simple. A tallying problem  $\exists \sigma s \leq t$  is turned into finding all substitutions  $\sigma$  for which  $s \setminus t$  is empty. The tallying procedure recursively explores the DNF of  $s \setminus t$ , generating a set of sets of constraints. When performed naively, this procedure may be suboptimal for two reasons. First, it may keep *unsatisfiable* constraints during constraint generation to later on discard them. Second, it may keep *redundant* constraints which yield correct, but useless substitutions. For instance, in the set of substitutions  $\{\{\alpha \mapsto \text{int}\}, \{\alpha \mapsto \mathbb{1}\}\}$ , the second one is useless (although correct) since the first one is more precise. As our evaluation shows (Section 7), generating such sets of sets of constraints naively is very costly to the point that some simple tallying problems cannot be solved in reasonable time. We therefore propose an alternative presentation which incorporates eager simplifications of constraints and greatly reduces the number of returned substitutions. We first introduce *normalized constraint sets*.

*Definition 5.1 (Normalized constraint set).* A normalized constraint set  $C$  is a set of triples  $\{(s_i, \alpha_i, t_i)\}_{i \in I}$  where all  $\alpha_i$  are distinct. The set  $\{\alpha_i\}_{i \in I}$  is called the *domain* of  $C$  and is written  $\text{dom}(C)$ . The empty constraint set is noted  $\top$ . Lastly, we define the constraint associated with a variable  $\alpha$  in a constraint  $C$  by:

$$C(\alpha) = (s, t) \text{ if } (s, \alpha, t) \in C \qquad C(\alpha) = (\mathbb{0}, \mathbb{1}) \text{ otherwise}$$

We now define operations on normalized constraint sets. These are parameterized by a set  $\Delta$  of type variables the tallying algorithm is allowed to instantiate.

*Definition 5.2 (Intersection of normalized constraint sets).* Let  $C_1$  and  $C_2$  be two normalized constraint sets and  $\Delta$  a set of type variables. We define the *intersection*  $C_1 \sqcap^\Delta C_2$  as follows. Let  $C = \{(s_1 \vee s_2, \alpha, t_1 \wedge t_2) \mid \alpha \in \text{dom}(C_1) \cup \text{dom}(C_2), (s_i, t_i) = C_i(\alpha) \text{ for } i = 1..2\}$ .

$$\begin{aligned} C_1 \sqcap^\Delta C_2 &= C & \text{if } \forall (s, \alpha, t) \in C, \text{vars}(s) \cup \text{vars}(t) \subseteq \Delta \Rightarrow s \leq t \\ C_1 \sqcap^\Delta C_2 &= \bullet & \text{otherwise} \end{aligned}$$

Intuitively, the intersection of two constraint sets is only defined if it does not create a trivially unsatisfiable constraint such as  $(\text{int}, \alpha, \text{bool})$ , which can never be satisfied.

We now want to define intersections and unions over sets of constraint sets, ranged over by  $\mathcal{C}$ —each constraint set  $C \in \mathcal{C}$  representing a potential solution of our tallying instance. In order to avoid a combinatorial explosion and to generate a minimal set of solutions, we want to eliminate redundant constraint sets from  $\mathcal{C}$ : in particular, there should not be two different constraint sets  $C_1, C_2 \in \mathcal{C}$  such that  $C_1$  *subsumes*  $C_2$ .

A first attempt at defining constraint set subsumption is to compare constraints for each variable independently, yielding the following definition:

*Definition 5.3 (Subsumption of constraint sets).* Let  $C_1$  and  $C_2$  be two normalized constraint sets. We say that  $C_1$  *subsumes*  $C_2$  (intuitively,  $C_1$  is more restrictive than  $C_2$ ), written  $C_1 \sqsubseteq C_2$ , if and only if:

$$\forall (s, \alpha, t) \in C_2. s \leq s' \text{ and } t' \leq t \text{ where } (s', t') = C_1(\alpha)$$

In other words, a set of constraints  $C_1$  *subsumes* a set of constraints  $C_2$  if  $C_1$  gives better bounds (higher lower bounds and lower upper bounds) for all variables of  $C_2$ . However, this definition of subsumption is not sufficient to avoid generating redundant solutions. Consider for instance the following two constraint sets:

$$C_1 = \{(\text{int}, \alpha, \mathbb{1}), (\mathbb{1}, \beta, \mathbb{1})\} \not\sqsubseteq \{(\text{int}, \alpha, \beta), (\text{int}, \beta, \mathbb{1})\} = C_2$$

These two constraint sets are not related by  $\sqsubseteq$ :  $C_2(\alpha)$  strictly subsumes  $C_1(\alpha)$ , while  $C_1(\beta)$  strictly subsumes  $C_2(\beta)$ . However, we should have  $C_1 \sqsubseteq C_2$ : while the constraint  $(\text{int}, \alpha, \mathbb{1})$  does not subsume  $(\text{int}, \alpha, \beta)$  by itself, it does when we consider it in a context where  $\beta \approx \mathbb{1}$  (which is enforced by the constraint  $(\mathbb{1}, \beta, \mathbb{1}) \in C_1$ ). To account for these indirect relations, we refine our definition of constraint set subsumption:

*Definition 5.4 (Subsumption of constraint sets, improved).* Let  $C_1$  and  $C_2$  be two normalized constraint sets. We say that  $C_1$  subsumes  $C_2$ , written  $C_1 \sqsubseteq C_2$ , if and only if:

$$\forall (s, \alpha, t) \in C_2. s \leq [s']_{C_1} \text{ and } [t']_{C_1} \leq t \text{ where } (s', t') = C_1(\alpha)$$

with  $[t]_C$  and  $[t]_C$  respectively a weakening and strengthening operator, inductively defined as follows:

$$\begin{aligned} [t]_{\{(s, \alpha, t)\} \cup C} &= [t\{\alpha \overset{+}{\rightsquigarrow} s\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \wedge t\}]_C && \text{where } \forall \alpha' \in \text{dom}(C). \alpha \preceq \alpha' \\ [t]_{\{(s, \alpha, t)\} \cup C} &= [t\{\alpha \overset{+}{\rightsquigarrow} s\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \vee s\}]_C && \text{where } \forall \alpha' \in \text{dom}(C). \alpha \preceq \alpha' \end{aligned}$$

where  $t\{\alpha \overset{+}{\rightsquigarrow} s\}$  denotes the type  $t$  where every positive top-level occurrence of  $\alpha$  has been substituted by  $s$ , and  $t\{\alpha \overset{-}{\rightsquigarrow} s\}$  denotes the type  $t$  where every negative top-level occurrence of  $\alpha$  has been substituted by  $s$  (both operations can easily be implemented on the BDT representation of the polymorphic descriptor of  $t$ ).

The side condition  $\forall \alpha' \in \text{dom}(C). \alpha \preceq \alpha'$  ensures that constraints are propagated following their dependency order: the tallying algorithm preserves the invariant that any constraint  $(s, \alpha, t)$  is such that  $\text{top\_vars}(s)$  and  $\text{top\_vars}(t)$  only contain type variables that are greater than  $\alpha$  by  $\preceq$ . See Appendix E for a formal proof that  $C_1 \sqsubseteq C_2$  implies every solution of  $C_1$  satisfies  $C_2$ .

As an example, consider our two constraint sets from earlier,  $C_1 = \{(\text{int}, \alpha, \mathbb{1}), (\mathbb{1}, \beta, \mathbb{1})\}$  and  $C_2 = \{(\text{int}, \alpha, \beta), (\text{int}, \beta, \mathbb{1})\}$ . This time, we do have  $C_1 \sqsubseteq C_2$ : when checking the constraints on  $\alpha$ , we are comparing  $(\text{int}, \alpha, \mathbb{1})$  against  $([\text{int}]_{C_1}, \alpha, [\beta]_{C_1}) = (\text{int}, \alpha, \mathbb{1})$ . We are now equipped to define a simplification function over sets of normalized constraint sets:

*Definition 5.5 (Simplification of sets of constraint sets).* The simplification of a set  $\mathcal{C}$  of constraint sets, written  $\text{csimp}(\mathcal{C})$ , is defined as follows:

$$\begin{aligned} \text{csimp}(\mathcal{C}) &= \text{csimp}'(\mathcal{C}, \emptyset) \\ \text{csimp}'(\emptyset, \mathcal{D}) &= \mathcal{D} \\ \text{csimp}'(\{C\} \cup \mathcal{C}, \mathcal{D}) &= \text{csimp}'(\mathcal{C}, \mathcal{D}) && \text{if } \exists C' \in \mathcal{D}, C' \sqsubseteq C \\ \text{csimp}'(\{C\} \cup \mathcal{C}, \mathcal{D}) &= \text{csimp}'(\mathcal{C}, \{C\} \cup \mathcal{D} \setminus \mathcal{D}') && \text{where } \mathcal{D}' = \{C' \mid C' \in \mathcal{D}, C \sqsubseteq C'\} \end{aligned}$$

Simplified sets of constraint sets ensure that their elements are pairwise non-subsumable, each element denotes a different, incomparable substitution. Maintaining simplified sets is a key ingredient to tame the complexity of the tallying operation.

*Definition 5.6 (Union and intersection of sets of constraint sets).* Given two sets of normalized constraint sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$  and a set of type variables  $\Delta$ , we define their intersection and union as:

$$\begin{aligned} \mathcal{C}_1 \sqcap^\Delta \mathcal{C}_2 &= \text{csimp}(\{C_1 \sqcap^\Delta C_2 \mid C_1 \in \mathcal{C}_1, C_2 \in \mathcal{C}_2, C_1 \sqcap^\Delta C_2 \neq \bullet\}) && \text{(intersection)} \\ \mathcal{C}_1 \sqcup \mathcal{C}_2 &= \text{csimp}(\mathcal{C}_1 \cup \mathcal{C}_2) && \text{(union)} \end{aligned}$$

We can now review the tallying algorithm, Algorithm 3, which consists of three steps. The function `normalize` transforms a subtyping constraint  $N_1 \leq N_2$  into a set of normalized constraint sets  $\mathcal{C}$ . Like in the subtyping algorithm, the initial subtyping constraint is encoded as an emptiness constraint for the type  $N_1 \setminus N_2$ .

ALGORITHM 3 (TALLYING).

```

1 let tally  $\Delta$   $S =$ 
2   let  $\mathcal{C} = \prod^{\Delta}_{(N_1, N_2) \in S}$  normalize  $\Delta$   $(N_1 \setminus N_2)$  in
3   let  $\mathcal{C}' = \prod_{C \in \mathcal{C}}$  propagate  $\Delta \top C \emptyset$  in
4   { solve  $C \mid C \in \mathcal{C}'$  }

```

Each normalized constraint set is then processed through the propagate function, which eliminates unsatisfiable constraint sets and recursively enriches satisfiable ones with new induced constraints: for every constraint  $(s, \alpha, t)$  produced by normalize, if either  $s$  or  $t$  contains type variables that are not in  $\Delta$ , then  $s \leq t$  may require new constraints. These new constraints are added by propagate to form a final set of constraint sets  $\mathcal{C}'$ . Each constraint set in  $\mathcal{C}'$  can then be turned into a set of recursive type equations and solved,

yielding a substitution. The code of normalize and propagate is given in Figure 5 (the code of solve is straightforward).

*Normalize.* The implementation of `normalize  $\Delta$   $N$`  follows the same structure as the emptiness test. It uses a hash table  $M$  to map already visited nodes to their set of constraint sets (l. 3). It plays a similar role as  $\Sigma$  in the subtyping algorithm. Whenever the emptiness test (which computes a Boolean) uses an “and” (resp. an “or”), the normalize function uses  $\sqcap$  (resp.  $\sqcup$ ). The first time a type

```

1 let  $M = \text{H.create } ()$ 
2 let rec normalize  $\Delta$   $N =$ 
3   if  $N \in M$  then H.find  $M$   $N$ 
4   else begin
5     H.add  $M$   $N$  { $\top$ };
6     let  $\mathcal{C} = \prod^{\Delta}_{(A_p, A_n, D) \in \text{dnf}(N.\text{def})}$  normalize_summand  $\Delta$   $A_p$   $A_n$   $D$  in
7     H.remove  $M$   $N$ ;  $\mathcal{C}$  end
8 and normalize_summand  $\Delta$   $A_p$   $A_n$   $D =$ 
9   let  $A = (A_p \cup A_n) \setminus \Delta$  in
10  if  $A \neq \emptyset$  then
11    let  $\alpha = \min A$  in (* Minimal type variable according to  $\leq$  *)
12    if  $\alpha \in A_p$  then {{ $(\emptyset, \alpha, \neg(A_p \setminus \{\alpha\} \wedge A_n \wedge D))$ }} else {{ $(A_p \wedge A_n \setminus \{\alpha\} \wedge D, \alpha, \mathbb{1})$ }}
13    else (norm_const  $\Delta$   $D.\text{const}$ )  $\sqcap^{\Delta}$  (norm_prod  $\Delta$   $D.\text{prod}$ )  $\sqcap^{\Delta}$  ...
14 let rec propagate  $\Delta$   $C_0$   $C_1$   $\Sigma =$ 
15  if  $C_1 = \top$  then { $C_0$ } else
16    let  $(N_1, \alpha, N_2) :: C'_1 = C_1$  in
17    let  $N = N_1 \setminus N_2$  in
18    if  $N \in \Sigma$  then propagate  $\Delta$  ( $C_0 \sqcap^{\Delta} \{(N_1, \alpha, N_2)\})$   $C'_1$   $\Sigma$ 
19    else
20      let  $\mathcal{C} = \{C_0 \sqcap C_1\} \sqcap^{\Delta}$  (normalize  $\Delta$   $N$ ) in
21       $\prod_{C \in \mathcal{C}}$  propagate  $\Delta \top C$  ( $\Sigma \cup \{N\}$ )

```

Fig. 5. The normalize and propagate functions

is visited (l. 5), the satisfiable constraint is added to the memoization table  $M$ . Each conjunction of the DNF of the type is then explored to generate individual sets of constraint sets which are intersected (l. 6). The node is removed from the table  $M$  and the resulting constraints returned.

The emptiness of a conjunction of the form  $\bigwedge_{\alpha \in A_p} \alpha \wedge \bigwedge_{\alpha \in A_n} \neg \alpha \wedge D$  is expressed by the function `normalize_summand`. The function chooses the smallest (according to the total order  $\preceq$  on type variables) top-level variable that is not in  $\Delta$ , and extract it from the summand, yielding a single constraint on that variable denoting the emptiness of the summand (l. 10–12). If there are no top-level variable, the function recursively explores the underlying constructors, with a function that mimics the emptiness test on components (and recursively calls `normalize`).

*Propagate.* The function `propagate`  $\Delta C_0 C_1 M \Sigma$  iteratively enriches the constraint set  $C_0 \sqcap C_1$  with the constraints induced by  $C_1$  ( $C_0$  contains the constraints already treated) until it reaches a fixpoint ( $\Sigma$  stores the induced constraints already visited). Note the interplay between `normalize` and `propagate`. When `normalize` generates a single constraint for a top-level variable, `propagate` calls `normalize` again on the bounds of that constraint.

*Solve.* The function `solve`  $C$  represents the constraint set  $C$  as a principal type substitution. It makes use of the `build` function (Section 4) to solve a set of equations. Each equation is generated by turning a subtyping constraint  $N_1 \leq \alpha \leq N_2$  into an equation  $\alpha = (N_1 \vee \alpha') \wedge N_2$ , for a fresh  $\alpha'$ .

## 6 Extensions

In Section 4 and Section 5, we described three core components: arrows, products, and constants. These components are the most commonly used in set-theoretic type formalizations, but a practical implementation of set-theoretic types also need other core components to model common programming structures (Section 6.1). Additionally, these core components can be used as building blocks and combined to define new components via encodings (Section 6.2).

### 6.1 Core extensions

*Intervals.* The component for constants cannot be used to represent all the individual integers, as the type `int` would then need to be defined as an infinite (or very large) union. Therefore, SSTT features a component representing intervals as sorted lists of disjoint intervals.

*Tuples.* Products can be generalized to cover tuples of any arity. For that, we can implement a generic tuple component parameterized by an arity  $n$ . It is straightforward to extend the subtyping and tallying algorithms accordingly (e.g. for  $n = 3$ , a DNF of tuples  $(t_1, t_2, t_3)$  is empty if and only if the corresponding DNF of pairs  $t_1 \times (t_2 \times t_3)$  is empty).

*Records.* The theory of set-theoretic types features records that can be open or closed, and fields that may be optionally absent. Record atoms are simply maps from field names to optional types, as well as a Boolean indicating if the record is closed (its list of fields is exact) or opened (it may have other, unknown fields). In order to represent absent fields, we define optional types as the pairs  $(t, a)$  where  $a$  is either  $\top$  (the value may be absent) or  $\perp$  (the value is not absent). An optional type  $(t, a)$  is empty if and only if  $t$  is empty and  $a = \perp$ . Then, subtyping for records can be encoded as subtyping for tuples: the record atoms composing a DNF can be turned into tuple atoms of optional types by fixing an order over fields. For instance, the DNF  $\{\ell_1 : t_1\} \vee \{\ell_2 : ? t_2 \dots\}$  can be encoded as  $((t_1, \perp), (\emptyset, \top), (\emptyset, \top)) \vee ((\perp, \top), (t_2, \top), (\perp, \top))$  (the first component of the tuples represents the field  $\ell_1$ , the second component represents  $\ell_2$ , and the last component represents the other fields). Note that while our record can be opened (extensible via sub-typing) they do not feature row polymorphism, which requires substantial modifications of the tallying procedure [15].

*Tagged types.* A tagged type  $T(t)$  intuitively represents the values  $\{T(v) \mid v \in t\}$ , where  $T \in \text{Tags}$  is a unary constructor of our language disjoint from other values (in particular, two types  $T_1(t)$  and  $T_2(t)$  with  $T_1 \neq T_2$  are disjoint). Tagged types have several uses: (i) they can be used to represent boxed values, and in conjunction with unions, to implement algebraic data types where every constructor is disjoint from the others, (ii) they can be used as a building block to extend our type algebra with new type constructors in a modular way using encodings (cf. Section 6.2), and (iii) in addition to the representation of boxed values, they can be generalized to represent opaque data types with an invariant or covariant parameter, as described in [24]. Extending the subtyping and tallying algorithms is straightforward, following the inductive relations given in this last work.

## 6.2 Encodings

In addition to core extensions, new type constructors can be added by encoding them as combinations of core components. This approach has two advantages: (i) it allows the programmer to extend the type algebra in a modular way, without having to modify the core components of the library, and (ii) it does not require any addition to the subtyping and tallying algorithms. Designing such an extension requires to answer these three questions:

- (1) how to encode our new type constructors using the existing core components?
- (2) how to recognize the encoding in a type?
- (3) how to extract the parameters of our type constructors from their encoding?

As long as we are only interested in subtyping and tallying, only the first point matters. However, as soon as we want to be able to inspect our types (e.g. for pretty-printing, cf. Appendix B), it becomes necessary to recognize when a subtree of our type comes from an encoding, and how to destruct this encoding. In order to make our encodings easily recognizable in a type, each extension should be disjoint from the other values. For instance, encoding the list type  $[\alpha^*]$  as the type  $X = \text{nil} \vee \alpha \times X$  as we did in Section 1 is ambiguous, as it becomes unclear whether  $\text{int} \times \text{nil}$  should be printed as a pair or as a list. We solve this problem by using tagged type: each constructor of an extension is tagged by a distinct tag that ensures its disjointness w.r.t other types (e.g. list and pairs are kept apart) and can be used by the pretty-printer to defer printing to a specialized function, provided by the extension.

Extension	Constructor	Encoding
Characters	Single character 'a'	$\text{Chr}(i_a)$ where $i_a$ is the ASCII encoding of 'a'
	Char interval 'a'-'z'	$\text{Chr}(i_a \dots i_z)$
	Any (supertype)	$\text{Chr}(\text{int})$
Strings	Single string "abc"	$\text{Str}(c_{\text{abc}})$ where $c_{\text{abc}}$ is a constant representing "abc"
	Any (supertype)	$\text{Str}(\text{cst})$ where $\text{cst}$ is the supertype of constants
Lists (regexp)	Empty list []	$\text{Lst}(c_{\text{nil}})$ where $c_{\text{nil}}$ is a constant representing []
	Any (supertype)	$\mu X. \text{Lst}(\mathbb{1} \times X) \vee \text{Lst}(c_{\text{nil}})$
	$t$ followed by $l$	$\text{Lst}(t \times l)$
	$t^*$ followed by $l$	$\mu X. \text{Lst}(t \times X) \vee l$
	$[ t_1 ; t_2 ; t_3^* ]$	$\text{Lst}(t_1 \times \text{Lst}(t_2 \times (\mu X. \text{Lst}(t_3 \times X) \vee \text{Lst}(c_{\text{nil}}))))$

Fig. 6. Example of extensions and their encodings

As an example, Figure 6 proposes some extensions and the corresponding encodings. The list encoding allows us to support regular expression types purely as an extension. Although they are isomorphic to recursive products, we can extract from their DNF an automaton representation of the regular language they capture, which can then be decompiled into a regular expression using

Brzozowski Algebraic Method [5] before being pretty-printed to the user. Other extensions have also been implemented through encodings, such as Booleans, floating-point numeric types, and opaque data types supporting multiple parameters of different variances.

*Gradual types.* An important contribution of [23] is that gradual types (that may contain a dynamic component « ? ») can be fully represented by a pair of static types (its lower and upper materializations). A type system implementor can simply represent types as pair of SSTT types without any modification. This is the approach taken by [Elixir](#) [9].

## 7 Evaluation

The SSTT library is implemented in OCaml. It amounts to 6900 loc, of which 3700 are the core type algebra and subtyping (Sections 3, 4 and 5.1). The code for the tallying is 470 loc, the rest being extensions (Section 6) and the pretty-printer (Appendix B). The goal of this evaluation showcase the benefits of the optimizations we presented. Specifically, we demonstrate that

*Type representation:* representing types as BDT with semantic simplifications is the right choice; *Subtyping:* our optimized subtyping (Alg. 2) is a clear improvement w.r.t. the naive one (Alg. 1); *Tallying:* our implementation of tallying (Sec. 5.2) is a clear improvement w.r.t. [14].

We also compare SSTT to (to the best of our knowledge) the only known complete and efficient implementation of set-theoretic types, that is the one found in the CDuce compiler. We discuss the specific differences between CDuce and our implementation in the related work section.

### 7.1 Experiments

We used a Linux workstation running Ubuntu 25.04, an Intel i5-4690@3.50GHz CPU and 32 GB of ram. All the code (whether using SSTT or CDuce) is compiled in native mode, using OCaml 5.4.0. As we want to benchmark only the primitive operations offered by SSTT and *not* a particular type system implementation, we devised the following experiment. In a first phase we type-checked several program files with the MLSem prototype [24] that uses SSTT as a type library, and which is an improved version of [13] (which itself used CDuce as a backend for types). The version we ran was instrumented to record in JSON files all the tallying instances that were solved during type-checking. This was done for three corpuses (some examples are given in Appendix C) with the following features:

- A. *Hindley-Milner* (150 loc, 30 functions, 872 tallying instances) Generic functions in a Hindley-Milner style inspired from the List and Map modules of OCaml's standard library. Most functions are recursive and involve pattern-matching over list or balanced binary tree types, and their type is inferred. The tallying instances generated involve multiple type variables (due to parametric polymorphism and recursion), but few unions or intersections (unions are only used to encode ADTs with disjoint constructors).
- B. *Overloaded* (220 loc, 68 functions, 5514 tallying instances) Non-recursive and non-generic overloaded functions. For each, a type capturing its overloaded behavior is inferred using type narrowing techniques. The most overloaded function is an intersection of 26 arrows. The tallying instances involve many intersections and unions, but fewer type variables.
- C. *HM+Overloaded* (345 loc, 83 functions, 2358 tallying instances) Recursive, generic and overloaded functions. This corpus involves parametric polymorphism, ad-hoc polymorphism, and encodings to model mutable data-structures (references, arrays, etc.). The tallying instances generated involve unions, intersections, and multiple type variables. This corpus also contains about 800 tallying instances that do not have solutions (either because it contains purposely ill-typed functions, or because some patterns or type cases become unreachable for some domains of overloaded functions).

*First phase, tallying instance generation.* We tested four different flavors of type representation:

*BDT:* types consists of plain reduced BDT

*SS:* types are reduced BDT with semantic simplifications (Sec. 3.3)

*HC:* types are reduced and hash-consed BDT (syntactically equal subtrees are shared in memory and set-theoretic operations and subtyping cache their results in memoization tables, as is usual for BDD operations)

*SS+HC:* types are semantically simplified and hash-consed BDT

For all four configurations, we tried to type-check the three corpuses while recording the tallying problems solved during type-checking. The result of this first phase is that if types are not semantically simplified systematically (*SS* or *SS+HC*), inferred types (in particular arrows) get too large during type-checking, to a point where some tallying instances do not terminate after a minute. This phenomenon occurs for all three corpuses. Implementing hash-consing for all the components and layers of the node structure does not help, though it allows sharing about 50% of the nodes, thus reducing the size of some BDTs by making more atoms syntactically equivalent. Consequently, our benchmark JSON files used in the second phase have been generated by running the type-checker in a setting where types are always semantically simplified. The practical impact is that none of the types stored in JSON files contain redundant atoms.

*Second phase, tallying instance solving.* The result of replaying the JSON files in various modes is given in Table 1. Each tallying instance in the JSON file is first parsed into a simple AST (not timed). The “Building” line represents the cumulative time (in s) necessary to build the type representations from the AST for all instances in the file. The “Solving” line reports the cumulative time to perform tallying and apply the resulting substitution to the initial type. The “Total” line is the sum of the two previous times. The “# Sol.” line gives the total number of tallying solutions for the corpus. The “Size” line reports the cumulative size of types in the tallying instances and their results. More specifically, it is the sum of the size of all tallying problems, all substitution solutions and all result types (applications of the latter to the former) of the tested file, in bytes. The size of a type is obtained by serializing it with OCaml’s `Marshal` library. The key point is that sharing is preserved and therefore, the serialization gives a good approximation of the *unique* subtrees that constitute a type. The “Avg. Size” line is just the “Size” line averaged by the number of entries in the file. The “Peak Size” line reports the size in bytes of the largest type built during the process (including intermediary types built during tallying and discarded). It is unavailable for CDuce as it is not possible to obtain it without modifying the internals of the library. Lastly, the “Timeout” line reports the number of instances that could not be solved in under ten seconds. In such cases the reported timings only include the time for instances that did not timeout.

The table is composed of four parts. The first part, “Type representation” fixes the subtyping and tallying to their optimized version and changes the underlying type representation (for the same four variants described in the first phase). The second part, “Subtyping” uses the textbook subtyping algorithm (Alg. 1) for the two type representation of reduced BDT and semantically simplified (*SS*) BDT. The third part, “Tallying” disables the constraint simplifications of Section 5.2, thus generating the same number of constraints as in [14]. Lastly, the “CDuce” column is the one where SSTT is replaced by CDuce and uses CDuce’s types, subtyping and tallying. Some type constructors are not supported by CDuce and have been encoded (cf. Appendix C).

## 7.2 Analysis

*On mandatory semantic simplifications.* As the results of the tallying constraint generation indicate, some form of semantic simplification is needed when dealing with set-theoretic types. Without these simplifications, BDTs become full of redundant atoms, to a point that traversing them

Corpus	Measure	Type representation				Naive Subtyping		Naive Tallying	CDuce
		BDT	SS*	HC	SS+HC	BDT	SS	SS	
A. 872 instances	Building	0.028	0.085	0.078	0.283	0.027	0.104	0.085	0.042
	Solving	0.258	0.352	0.481	0.865	0.557	1.015	8.524	0.480
	Total	0.286	0.437	0.559	1.147	0.584	1.119	8.609	0.522
	# Sol.	876	876	876	876	876	876	5250	880
	Size	3.3M	3.1M	2.4M	2.2M	3.3M	3.0M	11.0M	3.4M
	Avg. Size	3.8k	3.5k	2.7k	2.5k	3.8k	3.5k	12.7k	3.9k
	Peak Size	18.9k	9.3k	11.0k	6.7k	18.9k	9.3k	9.3k	N/A
Timeout	0	0	0	0	0	1	2	0	
B. 5514 instances	Building	0.015	0.067	0.053	0.204	0.013	0.074	0.065	0.045
	Solving	0.153	0.246	0.390	0.580	0.150	0.256	0.792	0.395
	Total	0.168	0.314	0.443	0.784	0.163	0.330	0.857	0.440
	# Sol.	5784	5784	5784	5784	5784	5784	10804	5784
	Size	3.9M	3.7M	3.4M	3.1M	3.9M	3.7M	4.8M	2.8M
	Avg. Size	714	665	614	562	714	665	872	501
	Peak Size	8.3k	8.5k	3.8k	3.8k	8.3k	8.5k	8.5k	N/A
Timeout	0	0	0	0	0	0	5	0	
C. 2358 instances	Building	0.021	0.067	0.058	0.214	0.021	0.072	0.064	0.036
	Solving	0.218	0.316	0.411	0.721	0.242	0.508	15.867	0.576
	Total	0.239	0.383	0.469	0.934	0.263	0.580	15.932	0.612
	# Sol.	1582	1582	1582	1582	1582	1582	17596	1655
	Size	2.7M	2.4M	2.2M	1.9M	2.7M	2.4M	19.6M	3.6M
	Avg. Size	1.1k	998	929	813	1.1k	999	8.3k	1.5k
	Peak Size	20.7k	6.9k	13.1k	4.8k	20.7k	6.9k	6.9k	N/A
Timeout	0	0	0	0	0	0	5	0	

★ : Default mode of the SSTT library

Table 1. Time performance (in s) and size of types (in bytes) for different configurations

during tallying becomes too costly. Instrumenting the first phase we found that some intermediary types peak at 50MB if not simplified. Simple hash-consing does not help either. Although it allows to share some nodes, these nodes are still present in the BDT and are traversed during tallying.

*Type representations.* In the first group of columns of Table 1, the “BDT” column plays the role of an “ideal but unpractical” representation. Indeed, since we are in a situation where every input type was the serialization of a semantically simplified type, then the tallying instances tested do not contain redundant atoms. The building time then only consists of syntactically constructing the reduced BDT. Solving constraints (tallying) which involves calling the subtyping algorithm and reconstructing finitely many new types is also the fastest. Unsurprisingly, the internal size of types is larger in BDT mode but the size and average size do not differ much from the other type representations. One metric which is particularly bad, however is the “Peak” size of types for corpuses A. and C., where the peak size is respectively two and three times larger than for the SS representation. This can be explained by the fact that these corpuses feature polymorphic input and output types, where many type substitutions occur. These may accumulate during type-checking and create very large types (in particular arrow types) which then make tallying impractical (e.g. while type-checking function composition or curried function applications). On the other hand, with semantic simplification, some cost is paid upfront during type building (or when building intermediary types during tallying), but this pays off as it allows performing tallying reasonably efficiently while preventing types from growing too much. Lastly, regarding hash-consing (with or without semantic simplification), it seems that the overhead caused by the look-up in memoization

tables outweighs the benefits of the reduced type size. We therefore conclude that systematic semantic simplification is a good default strategy which balances type size with runtime efficiency.

*Subtyping.* The second group of columns shows the impact of our improved subtyping algorithm. First for tallying instances that directly translate to simple subtyping tests (corpus B), switching to the naive algorithm shows the opposite effect, the naive algorithm seems faster (in BDT mode). This is because types are so simple that the caching machinery simply adds overhead. However the benefit of caching is real and adds up. For instances where the subtyping procedure is called repetitively either by tallying or by the semantic simplification of type, reverting to the naive subtyping algorithm has a dramatic effect (a simplification test in corpus A. even times out).

*Tallying.* To test our improvement w.r.t. [14], we deactivated the simplification of constraints (we do not apply `csimp` when merging constraint sets). Here the results are clear: it is the only mode where several examples failed due to timeouts (all examples involve overloaded functions with 10 arrows or more). Also note the sharp increase in the “# Sol.” and “Size” lines. As soon as tallying computes non-trivial substitutions (corpuses A and C) the number of redundant solutions found by the naive tallying is 6 to 11 times larger than the reduced set found by our improved algorithm.

*CDuce.* Last but not least, our experiments show that our default configuration for SSTT outperforms CDuce’s default one. In CDuce, types are represented by *lazy BDDs* (see Section 8) whose performances are between plain BDT and SS. CDuce’s subtyping works by building an inhabitant of the type, is written in CPS style, and avoids backtracking by traversing the type in BFS rather than DFS. In hindsight, although it may explore a subtree only when necessary, it performs a lot of memory allocations of closures which are detrimental in practice. For the three corpuses, we are respectively 20%, 40% and 60% faster than CDuce. Last, CDuce also generates more constraints, although not as much as [14] as it performs an early and unpublished version of `csimp`.

## 8 Related work and future work

*Set-theoretic types in the CDuce language.* The [CDuce language](#) [2, 18] has been the first to implement monomorphic higher-order set-theoretic types (the earlier [XDuce](#) [21] only featured first-order, top-level functions and no arrow type), later extended with parametric polymorphism [14]. This implementation supports arrows, pairs, records, integer intervals, characters and XML types. The overall architecture of descriptors we use to represent set-theoretic types follows the one described in [6], though the latter does not feature type variables. More generally, CDuce’s internal type representation and subtyping algorithm were introduced with XML aware programming in mind (its internal type algebra features XML namespaces and XML document types). And although [14] introduced parametric polymorphism and tallying for CDuce, their implementation was never discussed nor evaluated. Our work provides this missing piece in the literature, providing efficient implementations of all the relevant algorithms and data-structures, in an open and modular library.

*Set-theoretic types in mainstream languages.* Two examples of set-theoretic type systems for mainstream languages are [Etylizer](#) [26] (for Erlang), and [Elixir](#), [9]. Both provide an implementation of set-theoretic types that supports arrows and a generalized form of records for typing maps. While they seem to rely on the same ideas that were already present in CDuce (some form of BDD to represent types for instance), making a direct comparison is still not possible, either because their implementation is too slow [26] or in the case of Elixir, because it does not yet support recursive types nor type variables.

*Lazy BDDs.* [CDuce](#) and [Elixir](#) use lazy BDDs for representing DNF. As explained in [6], lazy BDDs avoid a potential explosion in size when repeatedly applying unions, and this until an intersection,

difference or negation is performed. When it happens, the accumulated unions of the lazy BDD must be distributed to compute a plain BDD, the hope being that the intersection or difference still yields a simpler BDD by canceling some lazy unions. Our semantic simplification (Section 3.3) can be adapted to work on lazy BDDs, but doing so increases its complexity and voids Property 3.9 that allows to quickly compute simplified negations. Overall, even if both approaches aim at keeping type representations small, they have different trade-offs: our systematic simplification tends to increase the building time of BDDs but simplifies future manipulations. Lazy BDDs have linear time unions but an increased complexity for some operations (such as negations).

*Syntactic and algebraic approaches.* While there are clear *theoretical* differences between semantic subtyping and syntax based ones, the separation lies elsewhere when it comes to implementation. Either the type system enjoys only *disjoint* union types, as is the case of MLsub [17]. In such works, subtyping remains polynomial, since types are equivalent to deterministic automata, and subtyping amounts to inclusion of their language, which is polynomial. Or some form of non-disjoint union occurs, as is the case in the early work of Aiken [1] or more recently MLStruct [25]. In both cases, either an exponential blow-up occurs when computing type operations (e.g. negation) or when deciding subtyping, by traversing the potentially exponential DNF of the type as we do. Both works mention in passing that memoizing intermediate results in a mutable table helps in practice, but do not give any detail. We believe our techniques could be applied to these works.

*SMT-based approaches.* Using SMT solvers to decide subtyping has been attempted before. For instance [3] uses an SMT solver to type-check the *D*-minor language. More recently, [29] uses an SMT solver to type-check Typescript, a typed dialect of JavaScript which features polymorphism, higher-order and intersection. While these works show that subtyping can to some extent be encoded into an SMT problem, it poses several challenges. First, [3] does not feature higher-order nor parametric polymorphism, which greatly simplifies the encoding. Second, it is not clear whether such encodings would rely on undecidable theories (and therefore would yield an incomplete subtyping). Also, as in [29], unification often depends on some ad-hoc procedure not directly handled by the SMT solver. In that work, type-checking a splay tree example program (18 functions, 206LOC and no refinement types) takes about 6 seconds which seems to indicate that, when refinement types are not needed, then the tailored algorithms we propose are better suited.

*Data-flow analysis.* The knowledgeable reader will remark that the subtyping algorithm could be cast in the framework of data-flow analysis. It is, after all, the computation of a fixed point of some monotonic predicate over a graph (the type, whose vertices are the nodes). For instance, Kildall’s algorithm [22] could be used to compute the emptiness predicate. The major difference with our approach is that, to the best of our knowledge, all variants of Kildall’s algorithm require the full graph. In our case, this would entail calling the  $\Psi_{\text{prod,arrow}}$  functions eagerly, generating a potentially exponential type up-front to test its emptiness.

*Conclusion and future work.* We proposed a modular data-structure to represent recursive set-theoretic types supporting multiple type constructors, unions, intersections, negations, and type substitutions. This representation relies on Binary Decision Trees (BDTs), equipped with a semantic simplification procedure which eliminates redundancy that accumulates when manipulating types, and we provide practical improvements to the traditional algorithms for subtyping and tallying. The performance of SSTT is compared against CDuce, the historical implementation of set-theoretic types, showing a significant improvement. While SSTT is already usable, we plan to extend it to support row polymorphism for record [7], as well as extensions relevant to dynamic languages such as Python (such as object types).

## Data availability statement

All the data-structures and algorithms described in this paper have been implemented in the OCaml library SSTT, available in the supplementary material, together with an extended version of this paper that includes appendices. A WebAssembly build that offers a REPL can be tested directly in the web browser. The artifact contains instructions for replicating the benchmarks. In case of acceptance, it will be submitted to artifact evaluation.

## References

- [1] Alexander Aiken and Brian R. Murphy. 1991. Implementing regular tree expressions. In *Functional Programming Languages and Computer Architecture*, John Hughes (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 427–447.
- [2] Véronique Benzaken, Giuseppe Castagna, and Alain Frisch. 2003. CDuce: an XML-centric general-purpose language. In *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming (Uppsala, Sweden) (ICFP '03)*. Association for Computing Machinery, New York, NY, USA, 51–63. doi:10.1145/944705.944711
- [3] Gavin M. Bierman, Andrew D. Gordon, Catalin Hritcu, and David E. Langworthy. 2012. Semantic subtyping with an SMT solver. *J. Funct. Program.* 22, 1 (2012), 31–105. doi:10.1017/S0956796812000032
- [4] Bryant. 1986. Graph-Based Algorithms for Boolean Function Manipulation. *IEEE Trans. Comput.* C-35, 8 (1986), 677–691. doi:10.1109/TC.1986.1676819
- [5] Janusz A. Brzozowski. 1964. Derivatives of Regular Expressions. *J. ACM* 11, 4 (Oct. 1964), 481–494. doi:10.1145/321239.321249
- [6] Giuseppe Castagna. 2020. Covariance and Contravariance: a fresh look at an old issue (a primer in advanced type systems for learning functional programmers). *Logical Methods in Computer Science* 16, 1 (2020), 15:1–15:58. doi:10.23638/LMCS-16(1:15)2020
- [7] Giuseppe Castagna. 2023. Typing Records, Maps, and Structs. *Proc. ACM Program. Lang.* 7, ICFP, Article 196 (Aug. 2023), 44 pages. doi:10.1145/3607838
- [8] Giuseppe Castagna. 2024. *Programming with Union, Intersection, and Negation Types*. Springer International Publishing, Cham, 309–378. doi:10.1007/978-3-031-34518-0\_12
- [9] Giuseppe Castagna, Guillaume Duboc, and José Valim. 2024. The Design Principles of the Elixir Type System. *The Art, Science, and Engineering of Programming* 8, 2 (2024). Video of the presentation I gave at the Erlang Symposium 2023: <https://youtu.be/VYmo867YF6g>. Also available, Guillaume Duboc’s presentation at ElixirConf EU 2023: <https://www.youtube.com/watch?v=gJJH7a2J9O8>.
- [10] Giuseppe Castagna and Alain Frisch. 2005. A gentle introduction to semantic subtyping. In *Proceedings of the 7th ACM SIGPLAN International Conference on Principles and Practice of Declarative Programming (Lisbon, Portugal) (PPDP '05)*. Association for Computing Machinery, New York, NY, USA, 198–208. doi:10.1145/1069774.1069793
- [11] Giuseppe Castagna, Victor Lanvin, Tommaso Petrucciani, and Jeremy G. Siek. 2019. Gradual Typing: a New Perspective. *Proc. ACM Program. Lang.* 3, Article 16, POPL ’19: 46th ACM Symposium on Principles of Programming Languages (jan 2019).
- [12] Giuseppe Castagna, Mickaël Laurent, and Kim Nguyen. 2024. Polymorphic Type Inference for Dynamic Languages. *Proc. ACM Program. Lang.* 8, POPL, Article 40 (Jan. 2024), 32 pages. doi:10.1145/3632882
- [13] Giuseppe Castagna, Mickaël Laurent, and Kim Nguyen. 2024. *Prototype: Polymorphic Type Inference for Dynamic Languages*. <https://doi.org/10.5281/zenodo.10155221>
- [14] Giuseppe Castagna, Kim Nguyen, Zhiwu Xu, and Pietro Abate. 2015. Polymorphic functions with set-theoretic types. Part 2: local type inference and type reconstruction. In *Proceedings of the 42nd Annual ACM SIGPLAN Symposium on Principles of Programming Languages (POPL '15)*. 289–302. doi:10.1145/2676726.2676991
- [15] Giuseppe Castagna and Loïc Peyrot. 2025. Polymorphic Records for Dynamic Languages. *Proc. ACM Program. Lang.* 9, OOPSLA1, Article 132 (April 2025), 28 pages. doi:10.1145/3720497
- [16] Giuseppe Castagna and Zhiwu Xu. 2011. Set-theoretic foundation of parametric polymorphism and subtyping. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming (Tokyo, Japan) (ICFP '11)*. Association for Computing Machinery, New York, NY, USA, 94–106. doi:10.1145/2034773.2034788
- [17] Stephen Dolan and Alan Mycroft. 2017. Polymorphism, Subtyping, and Type Inference in MLsub. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (Paris, France) (POPL '17)*. Giuseppe Castagna and Andrew D. Gordon (Eds.). Association for Computing Machinery, New York, NY, USA, 60–72. doi:10.1145/3009837.3009882
- [18] Alain Frisch. 2004. *Théorie, conception et réalisation d'un langage de programmation adapté à XML*. Ph.D. Dissertation. Université Paris Diderot.

- [19] Alain Frisch, Giuseppe Castagna, and Véronique Benzaken. 2008. Semantic subtyping: dealing set-theoretically with function, union, intersection, and negation types. *Journal of the ACM* 55, 4 (Sept. 2008), 19:1–19:64. <http://doi.acm.org/10.1145/1391289.1391293>
- [20] Nils Gesbert, Pierre Genevès, and Nabil Layaïda. 2015. A logical approach to deciding semantic subtyping. *ACM Transactions on Programming Languages and Systems* 38, 1 (2015), 3. doi:10.1145/2812805
- [21] Haruo Hosoya and Benjamin C. Pierce. 2003. XDuce: A statically typed XML processing language. *ACM Trans. Internet Technol.* 3, 2 (May 2003), 117–148. doi:10.1145/767193.767195
- [22] Gary A. Kildall. 1973. A unified approach to global program optimization. In *Proceedings of the 1st Annual ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages* (Boston, Massachusetts) (POPL '73). Association for Computing Machinery, New York, NY, USA, 194–206. doi:10.1145/512927.512945
- [23] Victor Lanvin. 2021. *A semantic foundation for gradual set-theoretic types*. Theses. Université Paris Cité. <https://theses.hal.science/tel-03853222>
- [24] Mickaël Laurent and Jan Vitek. 2026. Type Inference for Functional and Imperative Dynamic Languages. (2026). [https://mlaurent.ovh/publications/stt\\_implem.pdf](https://mlaurent.ovh/publications/stt_implem.pdf)
- [25] Lionel Parreaux and Chun Yin Chau. 2022. MLstruct: Principal Type Inference in a Boolean Algebra of Structural Types. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 141 (oct 2022), 30 pages. doi:10.1145/3563304
- [26] Albert Schimpf, Stefan Wehr, and Annette Bieniusa. 2023. Set-theoretic Types for Erlang. In *Proceedings of the 34th Symposium on Implementation and Application of Functional Languages* (Copenhagen, Denmark) (IFL '22). Association for Computing Machinery, New York, NY, USA, Article 4, 14 pages. doi:10.1145/3587216.3587220
- [27] Jeremy Siek and Walid Taha. 2006. Gradual typing for functional languages. *Scheme and Functional Programming*.
- [28] Sam Tobin-Hochstadt and Matthias Felleisen. 2008. The design and implementation of Typed Scheme. *ACM SIGPLAN Notices* 43, 1 (2008), 395–406.
- [29] Panagiotis Vekris, Benjamin Cosman, and Ranjit Jhala. 2016. Refinement types for TypeScript. *SIGPLAN Not.* 51, 6 (June 2016), 310–325. doi:10.1145/2980983.2908110

## A Semantic Subtyping

This appendix summarizes the interpretation of set-theoretic types formalized in [20] and how it can be used to define semantic subtyping.

### A.1 Type interpretation

In order to define subtyping over these types, the idea is to interpret each ground type (i.e., a type that does not contain type variables) as a set of values of our language. Then, subtyping can be defined as set containment over the interpretation of types. Intuitively, each ground type is associated to the set of values having this type: for instance, the base type `true` is interpreted as the singleton containing the constant `true`, while the type `bool = true ∨ false` is interpreted as the set `{true, false}`.

However, this idea becomes subtler when dealing with arrow types. Although an arrow type intuitively corresponds to a function (i.e., a  $\lambda$ -abstraction), interpreting an arrow type as a set of  $\lambda$ -abstractions is problematic as it yields a circular reasoning: determining if a  $\lambda$ -abstraction is in the interpretation of a type requires to define a type system, which in turns needs the subtyping relation that we are trying to build. In order to break this circularity, the interpretation of types is not defined over values of our language but over a domain  $\mathcal{D}$  defined below. Note that this does not necessarily invalidate the “types as set of values” intuition, as it is discussed in Castagna and Frisch [10, Section 2.7].

In addition, we need to define an interpretation for all types, and not only ground ones (the interpretation domain  $\mathcal{D}$  should account for type variables). A simple model was proposed by [20]. We succinctly present it in this section. The reader may refer to [8, Section 3.3] for more details.

*Definition A.1 (Interpretation domain [20]).* The *interpretation domain*  $\mathcal{D}$  is the set of finite terms  $d$  produced inductively by the following grammar

$$\begin{aligned} d &::= c^L \mid (d, d)^L \mid \{(d, \partial), \dots, (d, \partial)\}^L \\ \partial &::= d \mid \Omega \end{aligned}$$

where  $c$  ranges over the set  $C$  of constants,  $L$  ranges over finite sets of type variables, and where  $\Omega$  is such that  $\Omega \notin \mathcal{D}$ .

The elements of  $\mathcal{D}$  correspond, intuitively, to (denotations of) the results of the evaluation of expressions, labeled by finite sets of type variables. In particular, in a higher-order language, the results of computations can be functions which, in this model, are represented by sets of finite relations of the form  $\{(d_1, \partial_1), \dots, (d_n, \partial_n)\}^L$ , where  $\Omega$  (which is not in  $\mathcal{D}$ ) can appear in second components to signify that the function fails (i.e., evaluation is stuck) on the corresponding input. This is implemented by using in the second projection the meta-variable  $\partial$  which ranges over  $\mathcal{D}_\Omega = \mathcal{D} \cup \{\Omega\}$  (we reserve  $d$  to range over  $\mathcal{D}$ , thus excluding  $\Omega$ ). This constant  $\Omega$  is used to ensure that  $\mathbb{1} \rightarrow \mathbb{1}$  is not a supertype of all function types: if we used  $d$  instead of  $\partial$ , then every well-typed function could be subsumed to  $\mathbb{1} \rightarrow \mathbb{1}$  and, therefore, every application could be given the type  $\mathbb{1}$ , independently of its argument as long as this argument is typeable (see Section 4.2 of [19] for details). The restriction to *finite* relations corresponds to the intuition that the denotational semantics of a function is given by the set of its finite approximations, where finiteness is a restriction necessary (for cardinality reasons) to give the semantics to higher-order functions. Finally, the sets of type variables that label the elements of the domain are used to interpret type variables: we interpret a type variable  $\alpha$  by the set of all elements that are labeled by  $\alpha$ , that is the set  $\{d \mid \alpha \in \text{tags}(d)\}$  (where we define  $\text{tags}(c^L) = \text{tags}((d, d')^L) = \text{tags}(\{(d_1, \partial_1), \dots, (d_n, \partial_n)\}^L) = L$ ).

We now define the interpretation  $\llbracket t \rrbracket$  of a type  $t$  (this notation should not be mistaken with the interpretation  $\llbracket B \rrbracket$  of a BDD as a type, defined in Section 3). The interpretation  $\llbracket t \rrbracket$  should satisfy

the following equalities, where  $\mathcal{P}_{\text{fin}}$  denotes the restriction of the powerset to finite subsets, and  $\mathbb{B}$  denotes the function that assigns to each base type the set of constants of that type, so that for every constant  $c$  we have  $c \in \mathbb{B}(b_c)$  (we use  $b_c$  to denote the base type of the constant  $c$ ):

$$\begin{aligned} \llbracket 0 \rrbracket &= \emptyset & \llbracket \alpha \rrbracket &= \{d \mid \alpha \in \text{tags}(d)\} & \llbracket t_1 \vee t_2 \rrbracket &= \llbracket t_1 \rrbracket \cup \llbracket t_2 \rrbracket \\ \llbracket \neg t \rrbracket &= \mathcal{D} \setminus \llbracket t \rrbracket & \llbracket b \rrbracket &= \{c^L \mid c \in \mathbb{B}(b)\} & \llbracket t_1 \times t_2 \rrbracket &= \{(d_1, d_2)^L \mid d_1 \in \llbracket t_1 \rrbracket, d_2 \in \llbracket t_2 \rrbracket\} \\ \llbracket t_1 \rightarrow t_2 \rrbracket &= \{R^L \mid R \in \mathcal{P}_{\text{fin}}(\mathcal{D} \times \mathcal{D}_\Omega) \text{ s.t. } \forall (d, \partial) \in R. d \in \llbracket t_1 \rrbracket \implies \partial \in \llbracket t_2 \rrbracket\} \end{aligned}$$

Note that, even though we included  $\mathbb{1}$  and the intersection  $\wedge$  in the syntax of our types (Definition 2.1), those two can be defined from the other constructors:  $\mathbb{1} = \neg 0$  and  $t_1 \wedge t_2 = \neg(\neg t_1 \vee \neg t_2)$  (De Morgan's law). It is easy to see that, with these definitions, we have  $\llbracket \mathbb{1} \rrbracket = \mathcal{D}$  and  $\llbracket t_1 \wedge t_2 \rrbracket = \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket$ . Thus, it is not necessary to define an interpretation for them.

We cannot take the equations above directly as an inductive definition of  $\llbracket \cdot \rrbracket$  because types are not defined inductively but coinductively. Notice however that the contractivity condition of Definition 2.1 ensures that the binary relation  $\triangleright \subseteq \mathcal{T} \times \mathcal{T}$  defined by  $t_1 \vee t_2 \triangleright t_i$ ,  $t_1 \wedge t_2 \triangleright t_i$ ,  $\neg t \triangleright t$  is Noetherian. This gives an induction principle<sup>4</sup> on  $\mathcal{T}$  that we use combined with structural induction on  $\mathcal{D}$  to give the following definition, which validates the equalities above.

*Definition A.2 (Set-theoretic interpretation of types).* We define a binary predicate  $(d : t)$  (“the element  $d$  belongs to the type  $t$ ”), where  $d \in \mathcal{D}$  and  $t \in \mathcal{T}$ , by induction on the pair  $(d, t)$  ordered lexicographically. The predicate is defined as follows:

$$\begin{aligned} (c^L : b) &= c \in \mathbb{B}(b) \\ (d : \alpha) &= \alpha \in \text{tags}(d) \\ ((d_1, d_2)^L : t_1 \times t_2) &= (d_1 : t_1) \text{ and } (d_2 : t_2) \\ (\{(d_i, \partial_i), \dots, (d_n, \partial_n)\}^L : t_1 \rightarrow t_2) &= \forall i \in [1..n]. \text{ if } (d_i : t_1) \text{ then } (\partial_i : t_2) \\ (d : t_1 \vee t_2) &= (d : t_1) \text{ or } (d : t_2) \\ (d : \neg t) &= \text{not } (d : t) \\ (\partial : t) &= \text{false} && \text{otherwise} \end{aligned}$$

We define the *set-theoretic interpretation*  $\llbracket \cdot \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})$  as  $\llbracket t \rrbracket = \{d \in \mathcal{D} \mid (d : t)\}$ .

## A.2 Semantic subtyping

Now that we have a set-theoretic interpretation of types, we can define the subtyping preorder and its associated equivalence relation as follows.

*Definition A.3 (Subtyping relation).* We define the *subtyping* relation  $\leq$  and the *subtyping equivalence* relation  $\simeq$  as  $t_1 \leq t_2 \stackrel{\text{def}}{\iff} \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket$  and  $t_1 \simeq t_2 \stackrel{\text{def}}{\iff} (t_1 \leq t_2) \text{ and } (t_2 \leq t_1)$ .

This subtyping relation is decidable and is sometimes referred to as *semantic subtyping*, as it is not defined on the syntax of the type but on its interpretation.

For any two type substitutions  $\phi_1$  and  $\phi_2$ , we write  $\phi_1 \simeq \phi_2$  the pointwise subtyping equivalence of  $\phi_1$  and  $\phi_2$ . An important property of the interpretation above is that subtyping is preserved by type substitutions:

$$\forall t_1, t_2, \phi. t_1 \leq t_2 \implies t_1 \phi \leq t_2 \phi$$

<sup>4</sup>In a nutshell, we can do proofs and give definitions by induction on the structure of unions and negations—and, thus, intersections—but arrows, products, and base types are the base cases for the induction.

## B Pretty printing of types

Section 4 defines an internal representation for types that captures their semantic properties. This semantic representation makes it possible to implement set-theoretic connectives ( $\wedge$ ,  $\vee$ ,  $\neg$ ) efficiently, going from an algebraic representation of a type into a more semantic representation. However, when printing a set-theoretic type, we need to go the other direction: going from our internal representation back into an algebraic one. We propose in this section a quick description of the challenges and techniques we implemented to convert types into an algebraic form.

### B.1 Printing of descriptors

For each component of a descriptor (arrow, product, constant), we can easily extract a DNF of atoms directly from their BDT representation. However finding a concise algebraic representation for unions of components requires some more work. The naive approach consisting in systematically printing the union of all components may generate overly complicated notations: for instance, the type  $\mathbb{1}$  is printed as `AnyArrow | AnyPair | AnyConst`, and the type  $t \stackrel{\text{def}}{=} \neg(\text{int} \times \text{int})$  is printed as `AnyArrow | (AnyPair \ (int, int)) | AnyConst`. Instead, we propose the following algorithm for printing a descriptor  $D$  in a more concise way:

- (1) Each component of  $D$  is converted into an algebraic representation, and these algebraic representations are regrouped inside a union. We simplify this union by eliminating empty clauses. This yields an algebraic representation  $P_1$  for  $D$ . For instance, for the descriptor of  $t$ , this yields  $P_1 \stackrel{\text{def}}{=} \text{AnyArrow} \mid (\text{AnyPair} \setminus (\text{int}, \text{int})) \mid \text{AnyConst}$ .
- (2) Independently, the complement of each component of  $D$  is converted into an algebraic representation, and these algebraic representations are regrouped inside a union. For instance, for the descriptor of  $t$ , this yields `EmptyArrow | (int, int) | EmptyConst`. Again, we simplify this union by eliminating empty clauses, yielding `(int, int)`. We obtain an algebraic representation  $N$  for the negation of  $D$ .
- (3) We negate the algebraic representation  $N$  and perform trivial simplifications (such as elimination of double negations), yielding a representation  $P_2$  for  $D$ . In our example, we obtain  $P_2 \stackrel{\text{def}}{=} \neg(\text{int}, \text{int})$ .
- (4) We compare the length of  $P_1$  and  $P_2$  (using any relevant metric, for instance the number of nodes in their AST), and return the smaller one (in our case,  $P_2$ ).

### B.2 Printing of (recursive) nodes

Our types are represented by graphs of nodes that may contain cycles (in the case of recursive types). Thus, when printing a descriptor, it is not enough to recursively call the printer on each node referenced by its atoms, as this may result in an infinite inlining of recursive types. Instead, we follow the following strategy for printing a node  $N$ :

- (1) For each node in the connected component  $\text{deps}(N)$  of  $N$ , we compute an algebraic form for the top-level structure of its definition (referenced nodes are kept as symbolic references). For instance, when printing the type  $\mu X. ((\alpha \wedge (\text{int} \rightarrow \text{bool})) \times X) \vee \text{nil}$  represented in Figure 2, this leaves us with:  
`N1 where N1=(N2,N1)|nil and N2='a&(N3->N4) and N3=int and N4=bool`
- (2) Then, we try to inline the definition of each node ( $N_1, N_2, N_3, N_4$  in our example), but only when doing so reduces the length of the whole algebraic representation. This way, recursive references will not be inlined. In our example, the nodes  $N_2, N_3, N_4$  are inlined, yielding the algebraic expression `N1 where N1=( 'a&(int->bool), N1)|nil`.

## C Benchmarks

The corpuses have been type-checked using [24], a type system based on the techniques of inference and type narrowing introduced in [12].

*Compatibility with CDuce.* The subtyping constraints generated by the type-checker involve some type constructors that are not supported by CDuce. Thus, when running the benchmarks with the CDuce back end, the following encodings are used:

**Tuples** CDuce does not support tuples of arbitrary arity (it only supports pairs and records). Thus, a tuple  $(t_1, \dots, t_n)$  is encoded as a closed record  $\{\text{card} : n ; \_1 : t_1 ; \dots ; \_n : t_n\}$ .

**Tagged types** CDuce does not support *tagged types* [24] (in our benchmarks, tagged types are used to define constructors of algebraic data types). Thus, a tagged type  $T(t)$  is encoded as a pair  $(\_T, t)$  where  $\_T$  is a fresh constant for the constructor  $T$ .

**Opaque invariant constructors** Our benchmark HM+Overloaded relies on *opaque type constructors* [24] to encode references and mutable arrays. In CDuce, we encode a type  $\text{ref}(t)$  as a pair  $(\_ref, (t \rightarrow \text{unit}, \text{unit} \rightarrow t))$  where  $\_ref$  is a fresh constant.

We now include some excerpts from the corpuses.

### C.1 A. Hindley-Milner corpus

This corpus has been type-checked in a configuration where the type inference algorithm does not try to infer overloaded types for functions, but single arrow types (as in HM systems). Here is an excerpt (the type inferred for each function is added as a comment):

```

1 let fixpoint = fun f ->
2   let delta = fun x ->
3     f ( fun v -> ( x x v ))
4   in delta delta
5 (* (('a -> 'b) -> x1) -> x1 where x1 = 'c & ('a -> 'b) *)
6
7 let length_stub length lst =
8   if lst is [] then 0 else (length (tl lst))+1
9 (* ('a & [ any* ] -> int) -> [] | any::('a & [ any* ]) -> int *)
10
11 let length = fixpoint length_stub
12 (* [ any* ] -> int *)

```

Example of tallying instance generated when type-checking this excerpt:

```

"vars": [ "'a", "'b", "'c", "'d", "'e" ],
"mono": [],
"constr": [
  [
    "'(('a -> 'b) -> 'c & ('a -> 'b)) -> 'c & ('a -> 'b)",
    "'(('d & lst(x1) -> int) -> lst(tuple0 | (any, 'd & lst(x1))) -> int) -> 'e
      where x1 = tuple0 | (any, lst(x1))"
  ]
]

```

## C.2 B. Overloaded corpus

This corpus has been type-checked in a configuration where the type inference algorithm tries to infer overloaded types for functions (except when the domain of the function is explicitly given). Here is an excerpt (the type inferred for each function is added as a comment):

```

1 let succ x =
2   match x with
3   | 0 -> 1 | 1 -> 2 | 2 -> 3 | 3 -> 4 | 4 -> 5
4   | 5 -> 6 | 6 -> 7 | 7 -> 8 | 8 -> 9 | 9 -> 10
5   | 10 -> 11 | 11 -> 12 | 12 -> 13 | 13 -> 14 | 14 -> 15
6   | 15 -> 16 | 16 -> 17 | 17 -> 18 | 18 -> 19 | 19 -> 20
7   | 20 -> 21 | 21 -> 22 | 22 -> 23 | 23 -> 24 | 24 -> 25
8   | _ -> 0
9   end
10 (* (0 -> 1) & (1 -> 2) & ... & (24 -> 25) & (~ (0..24) -> 0) *)
11 let test (y:(5..15)) = succ (succ (succ (succ (succ y))))
12 (* (5..15) -> (10..20) *)

```

Example of tallying instance generated when type-checking this excerpt:

```

"vars": [ "'a" ],
"mono": [],
"constr": [
  [
    "(1 -> 2) & (2 -> 3) & (3 -> 4) & (4 -> 5)
    & (5 -> 6) & (6 -> 7) & (7 -> 8) & (8 -> 9) & (9 -> 10)
    & (10 -> 11) & (11 -> 12) & (12 -> 13) & (13 -> 14) & (14 -> 15)
    & (15 -> 16) & (16 -> 17) & (17 -> 18) & (18 -> 19) & (19 -> 20)
    & (20 -> 21) & (21 -> 22) & (22 -> 23) & (23 -> 24) & (24 -> 25)
    & (~ (0..24) -> 0) & (0 -> 1)",
    "(5..15) -> 'a"
  ]
]

```

## C.3 C. HM+Overloaded corpus

This corpus has been type-checked in a configuration where the type inference algorithm performs type narrowing, and where mutable data structures (arrays, references) are encoded as opaque data types. Here is an excerpt (the type inferred for each function is added as a comment):

```

1 let filter_imp (f:('a -> bool) & ('b -> false)) (arr:array('a|'b)) =
2   let res = array () in
3   let mut i = 0 in
4   while i < (len arr) do
5     let e = arr[i] in
6     if f e do push res e end ;
7     i := i + 1
8   end ;
9   return res
10 (* ('b -> false)&('a -> bool) -> array('b|'a) -> array('a\'b | 'c) *)

```

Example of tallying instance generated when type-checking this excerpt:

```
"vars": [ "'a", "'b", "'c" ],  
"mono": [ "'a", "'b" ],  
"constr": [  
  [  
    "('b -> bool(false)) & ('a -> bool(true | false))", "'a \\ 'b -> 'c"  
  ]  
]
```

### D Proofs for Section 3.3

This appendix contains detailed proofs of the three propositions stated in Section 3.3.

We recall the simplification definition for reference. Given a BDT  $B = N(a?B^+ : B^-)$  and a type context  $t$ , let  $B^{+'} = \text{simpl}(t \wedge a, B^+)$ ,  $B^{-'}$  =  $\text{simpl}(t \wedge \neg a, B^-)$ , and  $B' = N(a?B^{+'} : B^{-'})$ . Then:

$$\text{simpl}(t, \mathbb{L}(l)) = \mathbb{L}(l) \quad \text{and} \quad \text{simpl}(t, N(a?B^+ : B^-)) = \begin{cases} B^{+'} & \text{if } t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \\ B^{-'} & \text{else if } t \wedge \llbracket B^{-'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \\ B' & \text{otherwise.} \end{cases}$$

All three proofs establish a stronger invariant for  $\text{simpl}(t, B)$  with an arbitrary context  $t$ , from which the corresponding proposition follows by taking  $t = \mathbb{1}$ .

#### Proposition 3.8 (Correctness)

PROOF. We prove by structural induction on  $B$  the following invariant:

$$\text{for any type } t, \quad t \wedge \llbracket \text{simpl}(t, B) \rrbracket \simeq t \wedge \llbracket B \rrbracket. \quad (1)$$

*Leaf case.*  $B = \mathbb{L}(l)$ . Then  $\text{simpl}(t, \mathbb{L}(l)) = \mathbb{L}(l)$ , so the claim is immediate.

*Node case.*  $B = N(a?B^+ : B^-)$ . Set  $B^{+'} = \text{simpl}(t \wedge a, B^+)$ ,  $B^{-'}$  =  $\text{simpl}(t \wedge \neg a, B^-)$ , and  $B' = N(a?B^{+'} : B^{-'})$ . The induction hypotheses give

$$(t \wedge a) \wedge \llbracket B^{+'} \rrbracket \simeq (t \wedge a) \wedge \llbracket B^+ \rrbracket, \quad (2)$$

$$(t \wedge \neg a) \wedge \llbracket B^{-'} \rrbracket \simeq (t \wedge \neg a) \wedge \llbracket B^- \rrbracket. \quad (3)$$

Expanding  $\llbracket B' \rrbracket$  by definition and distributing  $t$  over the disjunction:

$$\begin{aligned} t \wedge \llbracket B' \rrbracket &= t \wedge \llbracket ([a] \wedge \llbracket B^{+'} \rrbracket) \vee (\neg[a] \wedge \llbracket B^{-'} \rrbracket) \rrbracket \\ &= (t \wedge a \wedge \llbracket B^{+'} \rrbracket) \vee (t \wedge \neg a \wedge \llbracket B^{-'} \rrbracket) \\ &\simeq (t \wedge a \wedge \llbracket B^+ \rrbracket) \vee (t \wedge \neg a \wedge \llbracket B^- \rrbracket) \quad \text{by (2) and (3)} \\ &= t \wedge \llbracket B \rrbracket. \end{aligned}$$

Each branch of  $\text{simpl}(t, B)$  therefore satisfies the invariant:

- If  $t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket$ :  $\text{simpl}(t, B) = B^{+'}$ , and  $t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \simeq t \wedge \llbracket B \rrbracket$ .
- If  $t \wedge \llbracket B^{-'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket$ :  $\text{simpl}(t, B) = B^{-'}$ , and  $t \wedge \llbracket B^{-'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket \simeq t \wedge \llbracket B \rrbracket$ .
- Otherwise:  $\text{simpl}(t, B) = B'$ , and  $t \wedge \llbracket B' \rrbracket \simeq t \wedge \llbracket B \rrbracket$ .

Taking  $t = \mathbb{1}$ :  $\llbracket \text{simpl}(B) \rrbracket \simeq \mathbb{1} \wedge \llbracket \text{simpl}(B) \rrbracket \simeq \mathbb{1} \wedge \llbracket B \rrbracket \simeq \llbracket B \rrbracket$ .  $\square$

#### Proposition 3.9 (Negation)

PROOF. We first establish an auxiliary lemma.

**Lemma.** For any BDT  $B$ ,  $\llbracket \neg B \rrbracket \simeq \neg \llbracket B \rrbracket$ .

*Proof.* By induction on  $B$ . The leaf case is immediate:  $\llbracket \neg \mathbb{L}(l) \rrbracket = \llbracket \mathbb{L}(\neg l) \rrbracket = \neg \llbracket l \rrbracket$ . For the node case  $B = N(a?B^+ : B^-)$ , note that  $\llbracket a \rrbracket$  and  $\neg \llbracket a \rrbracket$  partition the universe. Restricting the desired equality to each part:

$$\llbracket a \rrbracket \wedge \neg \llbracket B \rrbracket = \llbracket a \rrbracket \wedge \neg \llbracket B^+ \rrbracket \simeq \llbracket a \rrbracket \wedge \llbracket \neg B^+ \rrbracket \quad (\text{IH on } B^+),$$

$$\neg \llbracket a \rrbracket \wedge \neg \llbracket B \rrbracket = \neg \llbracket a \rrbracket \wedge \neg \llbracket B^- \rrbracket \simeq \neg \llbracket a \rrbracket \wedge \llbracket \neg B^- \rrbracket \quad (\text{IH on } B^-).$$

Combining:  $\neg \llbracket B \rrbracket = (\llbracket a \rrbracket \wedge \llbracket \neg B^+ \rrbracket) \vee (\neg \llbracket a \rrbracket \wedge \llbracket \neg B^- \rrbracket) = \llbracket N(a?\neg B^+ : \neg B^-) \rrbracket = \llbracket \neg B \rrbracket$ .  $\square$

We now prove by structural induction on  $B$  the following invariant:

$$\text{for any type } t, \quad \neg \text{simpl}(t, B) = \text{simpl}(t, \neg B) \quad (\text{syntactic equality}). \quad (4)$$

*Leaf case.*  $\neg\text{simpl}(t, \mathsf{L}(l)) = \neg\mathsf{L}(l) = \mathsf{L}(\neg l) = \text{simpl}(t, \mathsf{L}(\neg l)) = \text{simpl}(t, \neg\mathsf{L}(l))$ .

*Node case.*  $B = \mathsf{N}(a?B^+ : B^-)$ . Set  $B^{+'} = \text{simpl}(t \wedge a, B^+)$ ,  $B^{-'} = \text{simpl}(t \wedge \neg a, B^-)$ ,  $B' = \mathsf{N}(a?B^{+'} : B^{-'})$ . Since  $\neg B = \mathsf{N}(a?\neg B^+ : \neg B^-)$ , the computation of  $\text{simpl}(t, \neg B)$  proceeds as follows:

- Positive subtree:  $\text{simpl}(t \wedge a, \neg B^+) = \neg\text{simpl}(t \wedge a, B^+) = \neg B^{+'}$  (IH on  $B^+$ ).
- Negative subtree:  $\text{simpl}(t \wedge \neg a, \neg B^-) = \neg\text{simpl}(t \wedge \neg a, B^-) = \neg B^{-'}$  (IH on  $B^-$ ).
- Combined node:  $\mathsf{N}(a?\neg B^{+'} : \neg B^{-'}) = \neg B'$ .

It remains to show that the branching conditions for  $\text{simpl}(t, \neg B)$  agree with those for  $\text{simpl}(t, B)$ . By the auxiliary lemma,  $\llbracket \neg B^{+'} \rrbracket = \neg \llbracket B^{+'} \rrbracket$  and  $\llbracket \neg B^{-'} \rrbracket = \neg \llbracket B^{-'} \rrbracket$ , so the first condition for  $\text{simpl}(t, \neg B)$  reads  $t \wedge \neg \llbracket B^{+'} \rrbracket \simeq t \wedge \neg \llbracket B' \rrbracket$ . Since complementation within  $t$  is order-reversing and involutive,

$$t \wedge X \simeq t \wedge Y \iff t \wedge \neg X \simeq t \wedge \neg Y,$$

this condition holds if and only if  $t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket$ , i.e., the first condition for  $\text{simpl}(t, B)$ . The second condition is analogous. The three branches are therefore taken in the same order:

- First condition holds:  $\text{simpl}(t, B) = B^{+'}$  and  $\text{simpl}(t, \neg B) = \neg B^{+'}$ , so  $\neg\text{simpl}(t, B) = \text{simpl}(t, \neg B)$ .
- Second condition holds:  $\text{simpl}(t, B) = B^{-'}$  and  $\text{simpl}(t, \neg B) = \neg B^{-'}$ , so  $\neg\text{simpl}(t, B) = \text{simpl}(t, \neg B)$ .
- Otherwise:  $\text{simpl}(t, B) = B'$  and  $\text{simpl}(t, \neg B) = \neg B'$ , so  $\neg\text{simpl}(t, B) = \text{simpl}(t, \neg B)$ .

Taking  $t = \mathbb{1}$  yields the proposition.  $\square$

### Proposition 3.10 (Emptiness)

PROOF. We prove by structural induction on  $B$  the following invariant:

$$\text{for any type } t \neq \mathbb{0}, \quad t \wedge \llbracket B \rrbracket \simeq \mathbb{0} \iff \text{simpl}(t, B) = \perp. \quad (5)$$

*Leaf case.*  $B = \mathsf{L}(l)$ , so  $\text{simpl}(t, \mathsf{L}(l)) = \mathsf{L}(l)$ . We consider the two leaves of the Boolean leaf lattice  $\mathcal{L} = \{\perp, \top\}$ :

- $l = \perp$ :  $t \wedge \llbracket \mathsf{L}(\perp) \rrbracket = t \wedge \mathbb{0} = \mathbb{0}$ , and  $\mathsf{L}(\perp) = \perp$ . Both sides hold.
- $l = \top$ :  $t \wedge \llbracket \mathsf{L}(\top) \rrbracket = t \wedge \mathbb{1} = t \neq \mathbb{0}$  (by hypothesis), and  $\mathsf{L}(\top) \neq \perp$ . Both sides fail.

*Node case.*  $B = \mathsf{N}(a?B^+ : B^-)$ . Set  $B^{+'} = \text{simpl}(t \wedge a, B^+)$ ,  $B^{-'} = \text{simpl}(t \wedge \neg a, B^-)$ ,  $B' = \mathsf{N}(a?B^{+'} : B^{-'})$ .

*Direction ( $\Leftarrow$ ).* If  $\text{simpl}(t, B) = \perp$ , then by Proposition 3.8,

$$t \wedge \llbracket B \rrbracket \simeq t \wedge \llbracket \text{simpl}(t, B) \rrbracket = t \wedge \llbracket \perp \rrbracket = t \wedge \mathbb{0} = \mathbb{0}.$$

*Direction ( $\Rightarrow$ ).* Suppose  $t \wedge \llbracket B \rrbracket \simeq \mathbb{0}$ , i.e.,

$$(t \wedge a \wedge \llbracket B^{+'} \rrbracket) \vee (t \wedge \neg a \wedge \llbracket B^{-'} \rrbracket) \simeq \mathbb{0}. \quad (6)$$

This implies  $(t \wedge a) \wedge \llbracket B^{+'} \rrbracket \simeq \mathbb{0}$  and  $(t \wedge \neg a) \wedge \llbracket B^{-'} \rrbracket \simeq \mathbb{0}$ . Since  $t \neq \mathbb{0}$  and  $t \simeq (t \wedge a) \vee (t \wedge \neg a)$ , at most one of  $t \wedge a$  and  $t \wedge \neg a$  is empty. We distinguish three sub-cases.

- **Both  $t \wedge a \neq \mathbb{0}$  and  $t \wedge \neg a \neq \mathbb{0}$ :** By the induction hypothesis (each context is non-empty):  $B^{+'} = \perp$  and  $B^{-'} = \perp$ . The reduction property  $(\mathsf{N}(a?B : B) \rightsquigarrow B)$  then gives  $B' = \perp$ . The first condition  $t \wedge \llbracket B^{+'} \rrbracket \simeq t \wedge \llbracket B' \rrbracket$  reduces to  $\mathbb{0} \simeq \mathbb{0}$ , so  $\text{simpl}(t, B) = B^{+'} = \perp$ .
- $t \wedge a \simeq \mathbb{0}$  (so  $t \wedge \neg a \neq \mathbb{0}$ ): From (6),  $(t \wedge \neg a) \wedge \llbracket B^{-'} \rrbracket \simeq \mathbb{0}$ . By the induction hypothesis,  $B^{-'} = \perp$ . Moreover,

$$t \wedge \llbracket B' \rrbracket = (t \wedge a \wedge \llbracket B^{+'} \rrbracket) \vee (t \wedge \neg a \wedge \llbracket B^{-'} \rrbracket) \simeq \mathbb{0} \vee \mathbb{0} = \mathbb{0}.$$

In particular  $t \wedge \llbracket B^{-'} \rrbracket = \mathbb{0} \simeq t \wedge \llbracket B' \rrbracket$ , so the second condition triggers and  $\text{simpl}(t, B) = B^{-'} = \perp$ .

- $t \wedge \neg a \simeq 0$  (so  $t \simeq t \wedge a \neq 0$ ): Symmetric to the previous case: by the induction hypothesis  $B^{+'} = \perp$ , and  $t \wedge \llbracket B^{+'} \rrbracket = 0 \simeq t \wedge \llbracket B' \rrbracket$ , so  $\text{simpl}(t, B) = B^{+'} = \perp$ .

Taking  $t = \mathbb{1}$  (which satisfies  $\mathbb{1} \neq 0$ ) yields the proposition. □

## E Tallying: Constraint Sets Subsumption Soundness

This appendix proves that constraint set subsumption (Definition 5.4) is *sound*: if  $C_1 \sqsubseteq C_2$ , then every substitution satisfying  $C_1$  also satisfies  $C_2$ . In other words, a more restrictive constraint set can always be substituted for a less restrictive one without losing any solution.

### Definitions

We first recall that a normalized constraint set  $C$  collects upper and lower bounds on type variables: a triple  $(s, \alpha, t) \in C$  requires that any solution  $\sigma$  instantiate  $\alpha$  between  $s$  and  $t$ , taking into account how  $\sigma$  instantiates the other variables present in  $s$  and  $t$ .

*Definition E.1 (Solution).* A substitution  $\sigma$  *satisfies* a normalized constraint set  $C$ , written  $\sigma \models C$ , if for every  $(s, \alpha, t) \in C$ ,

$$s\sigma \leq \alpha\sigma \leq t\sigma.$$

### Key lemma

The weakening and strengthening operators  $\lceil u \rceil_C$  and  $\lfloor u \rfloor_C$  (Definition 5.4) replace the top-level type variables of  $u$  by their bounds from  $C$  in a monotone (resp. anti-monotone) fashion. The crucial semantic property is that both operators are transparent under any solution of  $C$ .

LEMMA E.2 (SEMANTIC TRANSPARENCY). *For any type  $u$ , normalized constraint set  $C$ , and substitution  $\sigma$  with  $\sigma \models C$ :*

$$\lceil u \rceil_C \sigma = u\sigma \quad \text{and} \quad \lfloor u \rfloor_C \sigma = u\sigma.$$

PROOF. We prove both claims simultaneously by induction on  $C$ .

*Base case.*  $C = \emptyset$ . Then  $\lceil u \rceil_\emptyset = u = \lfloor u \rfloor_\emptyset$  by definition, so the claims are immediate.

*Inductive step.*  $C = \{(s, \alpha, t)\} \cup C'$  with  $\alpha \preceq \alpha'$  for all  $\alpha' \in \text{dom}(C')$ . By definition:

$$\lceil u \rceil_C = \lceil u\{\alpha \overset{+}{\rightsquigarrow} \alpha \vee s\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \wedge t\} \rceil_{C'},$$

$$\lfloor u \rfloor_C = \lfloor u\{\alpha \overset{+}{\rightsquigarrow} \alpha \wedge t\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \vee s\} \rfloor_{C'}.$$

Since  $\sigma \models C$  we have, in particular,  $s\sigma \leq \alpha\sigma \leq t\sigma$ . For any occurrence of  $\alpha$  in  $u$ , both substitution pairs evaluate to  $\alpha\sigma$  under  $\sigma$ :  $(\alpha \vee s)\sigma = \alpha\sigma \vee s\sigma = \alpha\sigma$  (since  $s\sigma \leq \alpha\sigma$ ) and  $(\alpha \wedge t)\sigma = \alpha\sigma \wedge t\sigma = \alpha\sigma$  (since  $\alpha\sigma \leq t\sigma$ ).

Therefore, letting  $u_1 = u\{\alpha \overset{+}{\rightsquigarrow} \alpha \vee s\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \wedge t\}$  and  $u_2 = u\{\alpha \overset{+}{\rightsquigarrow} \alpha \wedge t\}\{\alpha \overset{-}{\rightsquigarrow} \alpha \vee s\}$ , we have  $u_1\sigma = u\sigma$  and  $u_2\sigma = u\sigma$ .

Since  $C' \subseteq C$ , we have  $\sigma \models C'$ . The induction hypothesis for  $\lceil \cdot \rceil$  gives  $\lceil u_1 \rceil_{C'} \sigma = u_1\sigma = u\sigma$ , and the induction hypothesis for  $\lfloor \cdot \rfloor$  gives  $\lfloor u_2 \rfloor_{C'} \sigma = u_2\sigma = u\sigma$ .  $\square$

### Main theorem

THEOREM E.3 (SUBSUMPTION SOUNDNESS). *Let  $C_1$  and  $C_2$  be normalized constraint sets with  $C_1 \sqsubseteq C_2$  (Definition 5.4). For any substitution  $\sigma$  with  $\sigma \models C_1$ , we have  $\sigma \models C_2$ .*

PROOF. Take any  $(s, \alpha, t) \in C_2$ . Let  $(s', t') = C_1(\alpha)$ . We must show  $s\sigma \leq \alpha\sigma \leq t\sigma$ .

By definition of  $C_1(\alpha)$ :

- If  $\alpha \in \text{dom}(C_1)$ : the triple  $(s', \alpha, t') \in C_1$ , so  $\sigma \models C_1$  gives  $s'\sigma \leq \alpha\sigma \leq t'\sigma$ .
- If  $\alpha \notin \text{dom}(C_1)$ :  $(s', t') = (\emptyset, \mathbb{1})$  by convention, so  $s'\sigma = \emptyset \leq \alpha\sigma$  and  $\alpha\sigma \leq \mathbb{1} = t'\sigma$  trivially.

In both cases  $s'\sigma \leq \alpha\sigma \leq t'\sigma$ .

The subsumption condition  $C_1 \sqsubseteq C_2$  at  $(s, \alpha, t)$  gives:

$$s \leq \lceil s' \rceil_{C_1} \quad \text{and} \quad \lfloor t' \rfloor_{C_1} \leq t.$$

**Lower bound.**

$$s\sigma \leq \lceil s' \rceil_{C_1} \sigma = s'\sigma \leq \alpha\sigma.$$

The first inequality holds because applying a substitution preserves subtyping (Proposition 2.2, i.e.  $s \leq \lceil s' \rceil_{C_1}$  implies  $s\sigma \leq \lceil s' \rceil_{C_1} \sigma$ ). The equality is Lemma E.2. The last inequality is established above.

**Upper bound.**

$$\alpha\sigma \leq t'\sigma = \lfloor t' \rfloor_{C_1} \sigma \leq t\sigma.$$

The first inequality is established above. The equality is Lemma E.2. The last inequality holds because  $\lfloor t' \rfloor_{C_1} \leq t$  implies  $\lfloor t' \rfloor_{C_1} \sigma \leq t\sigma$  (Proposition 2.2).

Since  $s\sigma \leq \alpha\sigma \leq t\sigma$  for every  $(s, \alpha, t) \in C_2$ , we conclude  $\sigma \models C_2$ . □

The theorem justifies the use of  $\sqsubseteq$  to eliminate redundant constraint sets in csimp: if  $C_1 \sqsubseteq C_2$  then every substitution produced from  $C_1$  already satisfies  $C_2$ , so  $C_2$  contributes no new solutions and can safely be discarded.