

# Set-theoretic Foundation of Parametric Polymorphism and Subtyping

Giuseppe Castagna<sup>1</sup>    Zhiwu Xu<sup>1,2\*</sup>

<sup>1</sup>CNRS, Laboratoire Preuves, Programmes et Systèmes, Univ Paris Diderot, Sorbonne Paris Cité, Paris, France.

<sup>2</sup>State Key Laboratory of Computer Science, Institute of Software, Chinese Academy of Science, Beijing, China

**Abstract.** We define and study parametric polymorphism for a type system with recursive, product, union, intersection, negation, and function types. We first recall why the definition of such a system was considered hard—when not impossible—and then present the main ideas at the basis of our solution. In particular, we introduce the notion of “convexity” on which our solution is built up and discuss its connections with parametricity as defined by Reynolds to whose study our work sheds new light.

**Categories and Subject Descriptors** D.3.3 [Programming Languages]: Language Constructs and Features—Polymorphism; Data types and structure; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs—type structure

**General Terms** Theory, Languages

**Keywords** Types, subtyping, polymorphism, parametricity, XML

## 1. Introduction

The standard approach to defining subtyping is to define a collection of syntax-driven subtyping rules that relate types with similar syntactic structure. However the presence of logical operators, such as unions and intersections, makes a syntactic characterization difficult, which is why a semantic approach is used instead: type  $t$  is a subtype of type  $s$  if the set of values denoted by  $t$  is a subset of the set of values denoted by  $s$ . The goal of this study is to extend this approach to parametric types. Since such type systems are at the heart of functional languages manipulating XML data, our study directly applies to them. Parametric polymorphism has repeatedly been requested to and discussed in various working groups of standards (eg, RELAX NG [5] and XQuery [6]) since it would bring not only the well-known advantages already demonstrated in existing functional languages, but also new usages peculiar to XML. A typical example is SOAP [21] that provides XML “envelopes” to wrap generic content. Functions manipulating SOAP envelopes are thus working on polymorphically typed objects encapsulated in XML structures. Polymorphic higher-order functions are also

needed, as shown by our practice of Ocsigen [1], a framework to develop dynamic web sites where web-site paths (uri) are associated to functions that take the uri parameters—the so-called “query strings” [7]—and return a Xhtml page. The core of the dynamic part of Ocsigen system is the function `register_new_service` whose (moral) type is:

$$\forall X \triangleleft \text{Params}. (\text{Path} \times (X \rightarrow \text{Xhtml})) \rightarrow \text{unit}$$

That is, it is a function that registers the association of its two parameters: a path in the site hierarchy and a function that fed with a query string that matches the description  $X$  (Params being the XML type of all possible query strings) returns an Xhtml page. Unfortunately, this kind of polymorphism is not available and the current implementation of `register_new_service` must bypass the type system (of OCaml, OCamlDuce [9], and/or CDuce [2]), losing all the advantages of static verification.

So why despite all this interest and motivations does no satisfactory solution exist yet? The crux of the problem is that, despite several efforts (eg, [15, 20]), it was not known how to define—and *a fortiori* how to decide—the subtyping relation for XML types in the presence of type variables. Actually, a purely semantic subtyping approach to polymorphism was believed to be impossible.

In this paper we focus on this problem and—though, henceforth we will seldom mention XML—solve it. The solution—which is more broadly applicable than just to an XML processing setting—is based on some strong intuitions and a lot of technical work. We follow this dichotomy for our presentation and organize it in two parts: an informal description of the main ideas and intuitions underlying the approach and the formal description of the technical development. More precisely, in Section 2 we first examine why this problem is deemed unfeasible or unpractical and simple solutions do not work (§2.1-2.3). Then we present the intuition underlying our solution (§2.4) and outline, in an informal way, the main properties that make the definition of subtyping possible (§2.5) as well as the key technical details of the algorithm that decides it (§2.6). We conclude this first part by giving some examples of the subtyping relation (§2.7) and discussing related work (§2.8). In Section 3 we present the key steps of the formal treatment to support the claims of Section 2 in particular the soundness and completeness of the algorithm and the decidability of the subtyping relation. We conclude our presentation with Section 4 where we discuss at length the connections between parametricity and our solution, as well as the new perspectives of research that our work opens. It may seem odd that we focus on a subtyping relation without defining any calculus to use it. Defining the polymorphic subtyping relation and defining a polymorphic calculus are distinct problems (though the latter depends on the former), and there is no doubt that the former is the very harder one. Once the subtyping relation is defined, it can be immediately applied to simple polymorphic calculi (eg, an extension with higher-order functions of the

\* This author was partially supported by the National Natural Science Foundation of China under Grant No. 61070038 and by a joint training program by CNRS and the Chinese Academy of Science.

language in [15]) but the definition of calculi that fully exploit the expressiveness of these types and of algorithms that (even partially) infer type annotations are different and very difficult problems that we leave for future work.

## 2. The key ideas

### 2.1 Regular types

XML types are essentially regular tree languages: an XML type is the set of all XML documents that match the type. As such they can be encoded by product types (for concatenation), union types (for regex unions) and recursive types (for Kleene's star). To type higher order functions we need arrow types and we also need intersection and negation types since in the presence of arrows they can no longer be encoded. Therefore, studying polymorphism for XML types is equivalent to studying it for the types that are coinductively (for recursion) produced by the following grammar:

$$t ::= b \mid t \times t \mid t \rightarrow t \mid t \vee t \mid t \wedge t \mid \neg t \mid \emptyset \mid \mathbb{1}$$

where  $b$  ranges over basic types (eg, `Bool`, `Real`, `Int`, ...), and  $\emptyset$  and  $\mathbb{1}$  respectively denote the empty (ie, that contains no value) and top (ie, that contains all values) types. In other terms, types are nothing but a propositional logic (with standard logical connectives:  $\wedge$ ,  $\vee$ ,  $\neg$ ) whose atoms are  $\emptyset$ ,  $\mathbb{1}$ , basic, product, and arrow types. We use  $\mathcal{T}$  to denote the set of all types.

In order to preserve the semantics of XML types as sets of documents (but also to give programmers a very intuitive interpretation of types) it is advisable to interpret a type as the set of all values that have that type. Accordingly, `Int` is interpreted as the set that contains the values `0`, `1`, `-1`, `2`, ...; `Bool` is interpreted as the set the contains the values `true` and `false`; and so on. In particular, then, unions, intersections, and negations (ie, type connectives) must have a set-theoretic semantics. Formally, this corresponds to define an interpretation function from types to the sets of some domain  $\mathcal{D}$  (for simplicity the reader can think of  $\mathcal{D}$  as the set of all values of the language), that is  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})$ , such that (i)  $\llbracket t_1 \vee t_2 \rrbracket = \llbracket t_1 \rrbracket \cup \llbracket t_2 \rrbracket$ , (ii)  $\llbracket t_1 \wedge t_2 \rrbracket = \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket$ , (iii)  $\llbracket \neg t \rrbracket = \mathcal{D} \setminus \llbracket t \rrbracket$ , (iv)  $\llbracket \emptyset \rrbracket = \emptyset$ , and (v)  $\llbracket \mathbb{1} \rrbracket = \mathcal{D}$ . Once we have defined such an interpretation, then the subtyping relation is naturally defined in terms of it:

$$t_1 \leq t_2 \stackrel{\text{def}}{\iff} \llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket$$

which, restricted to XML types, corresponds to the usual interpretation of subtyping as language containment.

### 2.2 Higher-order functions

All definitions above run quite smoothly as long as basic and product types are the only atoms we consider (ie, the setting studied by Hosoya and Pierce [16]). But as soon as we add higher-order functions, that is, arrow types, the definitions above no longer work:

1. If we take as  $\mathcal{D}$  the set of all values, then this must include also  $\lambda$ -abstractions. Therefore, to define the semantic interpretation of types we need to define the type of  $\lambda$ -abstractions (in particular of the applications that may occur in their bodies) which needs the subtyping relation, which needs the semantic interpretation. We fall on a circularity.
2. If we take as  $\mathcal{D}$  some mathematical domain, then we must interpret  $t_1 \rightarrow t_2$  as the set of functions from  $\llbracket t_1 \rrbracket$  to  $\llbracket t_2 \rrbracket$ . For instance if we consider functions as binary relations, then  $\llbracket t_1 \rightarrow t_2 \rrbracket$  could be the set

$$\{ f \subseteq \mathcal{D}^2 \mid (d_1, d_2) \in f \text{ and } d_1 \in \llbracket t_1 \rrbracket \text{ implies } d_2 \in \llbracket t_2 \rrbracket \} \quad (1)$$

or, compactly,  $\mathcal{P}(\overline{\llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket \rrbracket})$ , where the  $\overline{S}$  denotes the complement of the set  $S$  within the appropriate universe (in words,

these are the sets of pairs in which it is *not* true that the first projection belongs to  $\llbracket t_1 \rrbracket$  and the second does not belong to  $\llbracket t_2 \rrbracket$ ). But here the problem is not circularity but cardinality, since this would require  $\mathcal{D}$  to contain  $\mathcal{P}(\mathcal{D}^2)$ , which is impossible.

The solution to *both* problems is given by the theory of semantic subtyping [10], and relies on the observation that in order to use types in a programming language we do not need to know what types are, but just how they are related (by subtyping). In other terms, we do not require the semantic interpretation to map arrow types into the set in (1), but just to map them into sets that induce the same subtyping relation as (1) does. Roughly speaking, this turns out to require that for all  $s_1, s_2, t_1, t_2$ , the function  $\llbracket \_ \rrbracket$  satisfies the property:

$$\llbracket s_1 \rightarrow s_2 \rrbracket \subseteq \llbracket t_1 \rightarrow t_2 \rrbracket \iff \mathcal{P}(\overline{\llbracket s_1 \rrbracket \times \llbracket s_2 \rrbracket}) \subseteq \mathcal{P}(\overline{\llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket}) \quad (2)$$

whatever the sets denoted by  $s_1 \rightarrow s_2$  and  $t_1 \rightarrow t_2$  are. Equation (2) above covers only the case in which we compare two single arrow types. But, of course, a similar restriction must be imposed also when comparing arbitrary Boolean combinations of arrows. Formally, this can be enforced as follows. Let  $\llbracket \_ \rrbracket$  be a mapping from  $\mathcal{T}$  to  $\mathcal{P}(\mathcal{D})$ , we define a new mapping  $\mathbb{E}_{\llbracket \_ \rrbracket}$  as follows (henceforth we omit the  $\llbracket \_ \rrbracket$  subscript from  $\mathbb{E}_{\llbracket \_ \rrbracket}$ ):

$$\begin{aligned} \mathbb{E}(\emptyset) &= \emptyset & \mathbb{E}(\mathbb{1}) &= \mathcal{D} \\ \mathbb{E}(\neg t) &= \mathcal{D} \setminus \mathbb{E}(t) & \mathbb{E}(b) &= \llbracket b \rrbracket \\ \mathbb{E}(t_1 \vee t_2) &= \mathbb{E}(t_1) \cup \mathbb{E}(t_2) & \mathbb{E}(t_1 \times t_2) &= \llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket \\ \mathbb{E}(t_1 \wedge t_2) &= \mathbb{E}(t_1) \cap \mathbb{E}(t_2) & \mathbb{E}(t_1 \rightarrow t_2) &= \mathcal{P}(\overline{\llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket}) \end{aligned}$$

Then  $\llbracket \_ \rrbracket$  and  $\mathcal{D}$  form a set-theoretic *model* of types if, besides the properties (i-v) for the type connectives we stated in §2.1, the function  $\llbracket \_ \rrbracket$  also satisfies the following property:

$$\llbracket t_1 \rrbracket \subseteq \llbracket t_2 \rrbracket \iff \mathbb{E}(t_1) \subseteq \mathbb{E}(t_2) \quad (3)$$

which clearly implies (2). All these definitions yield a subtyping relation with all the desired properties: type connectives (ie, unions, intersections, and negations) have a set-theoretic semantics, type constructors (ie, products and arrows) behave as set-theoretic products and function spaces, and (with some care in defining the language and its typing relation) a type can be interpreted as the set of values that have that type. All that remains to do is:

1. show that a model exists<sup>1</sup> (easy) and
2. show how to decide the subtyping relation (difficult).

Both points are solved in [10] and the resulting type system is at the core of the programming language `CDuce` [2].

### 2.3 The problem and a naive (wrong) solution

The problem we want to solve in this paper is how to extend the approach described above when we add type variables (in bold):

$$t ::= \alpha \mid b \mid t \times t \mid t \rightarrow t \mid t \vee t \mid t \wedge t \mid \neg t \mid \emptyset \mid \mathbb{1}$$

where  $\alpha$  ranges over a countable set of type variables  $\mathcal{V}$ . We did not include any explicit quantification for type variables: in this work (as well as, all works in the domain we are aware of, foremost [15, 20]) we focus on prenex parametric polymorphism where type quantification is meta-theoretic. Once more, the crux of the problem is how to *define* the subtyping relation between two types that contain type variables. Since we know how to subtype closed types (ie, types without variables), then a naive solution is to reuse this relation by considering all possible ground instances of types with variables. Let  $\sigma$  denote a *ground substitution*, that is

<sup>1</sup> We do not need to look for a particular model, since all models induce essentially the same subtyping relation: see [8] for details.

a substitution from type variables to closed types. Then according to our naive definition two types are in subtyping relation if so are their ground instances:

$$t_1 \leq t_2 \stackrel{\text{def}}{\iff} \forall \sigma. \llbracket t_1 \sigma \rrbracket \subseteq \llbracket t_2 \sigma \rrbracket \quad (4)$$

(provided that the domain of  $\sigma$  contains all the variables occurring in  $t_1$  and  $t_2$ ). This closely matches the syntactic intuition of subtyping for prenex polymorphism according to which the statement  $t_1 \leq t_2$  is to be intended as  $\forall \alpha_1 \dots \alpha_n (t_1 \leq t_2)$ , where  $\alpha_1 \dots \alpha_n$  are all the variables occurring in  $t_1$  or  $t_2$ . Clearly, the containment on the right hand side of (4) is a *necessary* condition for subtyping. Unfortunately, considering it also as *sufficient* and, thus, using (4) to define subtyping yields a subtyping relation that suffers too many problems to be useful.

The first obstacle is that, as conjectured by Hosoya in [15], if the subtyping relation defined by (4) is decidable (which is an open problem), then deciding it is at least as hard as the satisfiability problem for set constraint systems with *negative constraints*, which is NEXPTIME-complete and for which, so far, no practical algorithm is known.

But even if the subtyping relation defined by (4) were decidable and Hosoya’s conjecture wrong, definition (4) yields a subtyping relation that misses the intuitiveness of the relation on ground types. This can be shown by an example drawn from [15]. For the sake of the example, imagine that our system includes singleton types, that is types that contain just one value, for every value of the language, and consider the following subtyping statement:

$$t \times \alpha \leq (t \times \neg t) \vee (\alpha \times t) \quad (5)$$

where  $t$  is a closed type.

According to (4) the statement holds if and only if  $t \times s \leq (t \times \neg t) \vee (s \times t)$  holds for every closed type  $s$ . It is easy to see that the latter holds if and only if  $t$  is a singleton type. This follows from the set theoretic property that if  $S$  is a singleton, then for every set  $X$ , either  $S \subseteq X$  or  $X \subseteq \bar{S}$ . By using this property on the singleton type  $t$ , we deduce that for every ground substitution of  $\alpha$  either  $\alpha \leq \neg t$  (therefore  $t \times \alpha \leq t \times \neg t$ , whence (5) follows) or  $t \leq \alpha$  (therefore  $t \times \alpha = (t \times \alpha \setminus t) \vee (t \times t)$  and the latter is contained component-wise in  $(t \times \neg t) \vee (\alpha \times t)$ , whence (5) holds again). Vice versa, if  $t$  contains at least two values, then substituting  $\alpha$  by any singleton containing a value of  $t$  disproves the containment.

More generally, (5) holds if and only if  $t$  is an *indivisible* type, that is, a non-empty type whose only proper subtype is the empty type. Singleton types are just an example of indivisible types, but in the absence of singleton types, basic types that are pairwise disjoint are indivisible as well. Therefore, while the case of singleton types is evocative, the same problem occurs in a language with just the Int type, too.

Equation (5) is pivotal in our work. It gives us two reasons to think that the subtyping relation defined by (4) is unfit to be used in practice. First, it tells us that in such a system deciding subtyping is at least as difficult as deciding the indivisibility of a type. This is a very hard problem (see [3] for an instance of this problem in a simpler setting) that makes us believe more in the undecidability of the relation, than in its decidability. Second, and much worse, it completely breaks parametricity yielding a completely non-intuitive subtyping relation. Indeed notice that in the two types in (5) the type variable  $\alpha$  occurs on the right of a product in one type and on the left of a product in the other. The idea of parametricity is that a function cannot explore arguments whose type is a type variable, it can just discard them, pass them to another function or copy them into the result. Now if (4) holds it means that by a simple subsumption a function that is parametric in its second argument can be considered parametric in its first argument instead. Understanding the intuition underlying this subtyping relation for type variables

(where the same type variable may appear in unrelated positions in two related types) seems out of reach of even theoretically-oriented programmers. This is why a semantic approach for subtyping polymorphic types has been deemed unfeasible and discarded in favor of partial or syntactic solutions (see related works in §2.8).

## 2.4 Ideas for a solution

Although the problems we pointed out in [15] are substantial, they do not preclude a semantic approach to parametric polymorphism. Furthermore the shortcomings caused by the absence of this approach make the study well worth of trying. Here we show that—paraphrasing a famous article by John Reynolds [19]—subtyping of polymorphism *is* set-theoretic.

The conjecture that we have been following since we discovered the problem of [15], and that is at the basis of all this work, is that *the loss of parametricity is only due to the behavior of indivisible types*, all the rest works (more or less) smoothly. The crux of the problem is that for an indivisible type  $t$  the validity of the formula

$$t \leq \alpha \quad \text{or} \quad \alpha \leq \neg t \quad (6)$$

can *stutter* from one subformula to the other (according to the assignment of  $\alpha$ ) losing in this way the uniformity typical of parametricity. If we can give a *semantic* characterization of models in which *stuttering* is absent, we believed this would have yielded a subtyping relation that is (i) semantic, (ii) intuitive for the programmer,<sup>2</sup> and (iii) decidable. The problem with indivisible types is that they are either completely inside or completely outside any other type. What we need, then, is to make indivisible types “splitable”, so that type variables can range over strict subsets of any type, indivisible ones included. Since this is impossible at a syntactic level, we shall do it at a semantic level. First, we replace ground substitutions with semantic (set) assignments of type variables,  $\eta : \mathcal{V} \rightarrow \mathcal{P}(\mathcal{D})$ , and add to interpretation functions a semantic assignment as a further parameter (as is customary in denotational semantics):

$$\llbracket \cdot \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^{\mathcal{V}} \rightarrow \mathcal{P}(\mathcal{D})$$

Such an interpretation (actually, the pair  $(\llbracket \cdot \rrbracket, \mathcal{D})$ ) is then a *set-theoretic model* if and only if for all assignments  $\eta$  it satisfies the following conditions (in bold the condition that shows the role of the assignment parameter  $\eta$ ):

$$\begin{aligned} \llbracket \alpha \rrbracket \eta &= \eta(\alpha) & \llbracket \neg t \rrbracket \eta &= \mathcal{D} \setminus \llbracket t \rrbracket \eta \\ \llbracket 0 \rrbracket \eta &= \emptyset & \llbracket t_1 \vee t_2 \rrbracket \eta &= \llbracket t_1 \rrbracket \eta \cup \llbracket t_2 \rrbracket \eta \\ \llbracket 1 \rrbracket \eta &= \mathcal{D} & \llbracket t_1 \wedge t_2 \rrbracket \eta &= \llbracket t_1 \rrbracket \eta \cap \llbracket t_2 \rrbracket \eta \\ \llbracket t_1 \rrbracket \eta &\subseteq \llbracket t_2 \rrbracket \eta & \iff \mathbb{E}(t_1) \eta &\subseteq \mathbb{E}(t_2) \eta \end{aligned}$$

(where  $\mathbb{E}()$  is extended in the obvious way to cope with semantic assignments). Then the subtyping relation is defined as follows:

$$t_1 \leq t_2 \stackrel{\text{def}}{\iff} \forall \eta \in \mathcal{P}(\mathcal{D})^{\mathcal{V}}. \llbracket t_1 \rrbracket \eta \subseteq \llbracket t_2 \rrbracket \eta \quad (7)$$

In this setting, every type  $t$  that denotes a set of at least two elements of  $\mathcal{D}$  can be split by an assignment. That is, it is possible to define an assignment for which a type variable  $\alpha$  denotes a subset of  $\mathcal{D}$  that is neither completely inside nor completely outside the interpretation of  $t$ . Therefore for such a type  $t$ , neither equation (6) nor, *a fortiori*, equation (5) hold. It is then clear that the stuttering of (6) is absent in every set-theoretic model in which all non-empty types—indivisible types included—denote *infinite* subsets of  $\mathcal{D}$ . Infinite denotations for non-empty types look as a possible, though

<sup>2</sup>For instance, type variables can only be subsumed to themselves and according to whether they occur in a covariant or contravariant position, to  $\mathbb{1}$  and to unions in which they explicitly appear or to  $\mathbb{0}$  and intersections in which they explicitly appear, respectively. De Morgan’s laws can be used to reduce other cases to one of these.

specific, solution to the problem of indivisible types. But what we are looking for is not a particular solution. We are looking for a semantic characterization of the “uniformity” that characterizes parametricity, in order to define a subtyping relation that is, we repeat, semantic, intuitive, and decidable.

This characterization is provided by the property of *convexity*.

## 2.5 Convexity

A set theoretic model  $(\llbracket \cdot \rrbracket, \mathcal{D})$  is *convex* if and only if for every finite set of types  $t_1, \dots, t_n$  it satisfies the following property:

$$\begin{aligned} \forall \eta. (\llbracket t_1 \rrbracket \eta = \emptyset \text{ or } \dots \text{ or } \llbracket t_n \rrbracket \eta = \emptyset) \\ \iff \\ (\forall \eta. \llbracket t_1 \rrbracket \eta = \emptyset) \text{ or } \dots \text{ or } (\forall \eta. \llbracket t_n \rrbracket \eta = \emptyset) \end{aligned} \quad (8)$$

This property is the cornerstone of our approach. As such it deserves detailed comments. It states that, given any finite set of types, if every assignment makes some of these types empty, then it is so because there exists one particular type that is empty for all possible assignments.<sup>3</sup> Therefore convexity forces the interpretation function to behave uniformly on its zeros (*ie*, on types whose interpretation is the empty set). Now, the zeros of the interpretation function play a crucial role in the theory of semantic subtyping, since they completely characterize the subtyping relation. Indeed  $s \leq t \iff \llbracket s \rrbracket \subseteq \llbracket t \rrbracket \iff \llbracket s \rrbracket \cap \llbracket \bar{t} \rrbracket \subseteq \emptyset \iff \llbracket s \wedge \bar{t} \rrbracket = \emptyset$ . Consequently, checking whether  $s \leq t$  is equivalent to checking whether the type  $s \wedge \bar{t}$  is empty; likewise, equation (3) in §2.2 is equivalent to requiring that for all  $t$  it satisfies  $\llbracket t \rrbracket = \emptyset \iff \mathbb{E}(t) = \emptyset$ . We deduce that convexity forces the subtyping relation to have a uniform behavior and, ergo, rules out non-intuitive relations such as the one in (5). This is so because convexity prevents stuttering, insofar as in every convex model  $(\llbracket t \wedge \bar{\alpha} \rrbracket \eta = \emptyset \text{ or } \llbracket t \wedge \alpha \rrbracket \eta = \emptyset)$  holds for all assignments  $\eta$  if and only if  $t$  is empty.

Convexity is the property we seek. The resulting subtyping relation is semantically defined and preserves the set-theoretic semantics of type connectives (union, intersection, negation) and the containment behavior of set-theoretic interpretations of type constructors (set-theoretic products for product types and set-theoretic function spaces for arrow types). Furthermore, the subtyping relation is not only semantic but also intuitive. First, it excludes non-intuitive relations by imposing a uniform behavior distinctive of the parametricity *à la* Reynolds: as we discuss at length in the conclusion, we push the analogy much farther since we believe that parametricity and convexity are connected, despite the fact that the former is defined in terms of *transformations* of related terms while the latter deals only with (subtyping) relations. Second, it is very easy to explain the intuition of type variables to a programmer:

For what concerns subtyping, a type variable can be considered as a special new user-defined basic type that is unrelated to any other atom but  $\mathbb{0}$  and  $\mathbb{1}$  and itself.<sup>4</sup> Type variables are special because their intersection with any ground type may be non-empty, whatever this type is.

<sup>3</sup> We dubbed this property *convexity* after convex formulas: a formula is convex if whenever it entails a disjunction of formulas, then it entails one of them. The  $\Rightarrow$  direction of (8) (the other direction is trivial) states the convexity of assignments with respect to emptiness:  $\eta \in \mathcal{P}(\mathcal{D})^\mathcal{V} \Rightarrow \bigvee_{i \in I} \llbracket t_i \rrbracket \eta = \emptyset$  implies that there exists  $h \in I$  such that  $\eta \in \mathcal{P}(\mathcal{D})^\mathcal{V} \Rightarrow \llbracket t_h \rrbracket \eta = \emptyset$ .

<sup>4</sup> This holds true even for languages with bounded quantification which, as it is well known, defines the subtyping relation for type variables. Bounded quantification does not require any modification to our system, since it can be encoded by intersections: a type variable  $\alpha$  bounded by a type  $t$  can be encoded by a fresh (unbounded) variable  $\beta$  by replacing  $\beta \wedge t$  for every occurrence of  $\alpha$ . We can do even more, since by using intersections we can impose different bounds to different occurrences of the same variable.

Of course, neither in the theoretical development nor in the subtyping algorithm type variables are dealt with as basic types. They need very subtle and elaborated techniques that form the core of our work. But this complexity is completely transparent to the programmer which can thus rely on a very simple intuition.

All that remains to do is (i) to prove the convexity property is not too restrictive, that is, that there exists at least one convex set-theoretic model and (ii) to show an algorithm that checks the subtyping relation. Contrary to the ground case, *both* problems are difficult. While their solutions require a lot of technical results (see Section 3), the intuition underlying them is relatively simple. For what concerns the existence of a convex set-theoretic model, the intuition can be grasped by considering just the logical fragment of our types, that is, the types in which  $\mathbb{0}$  and  $\mathbb{1}$  are the only atoms. This corresponds to the (classical) propositional logic where the two atoms represent, respectively, false and true. Next, consider the instance of the convexity property given for just two types,  $t_1$  and  $t_2$ . It is possible to prove that every non-degenerate Boolean algebra (*ie*, every Boolean algebra with more than two elements) satisfies it. Reasoning by induction it is possible to prove that convexity for  $n$  types is satisfied by any Boolean algebra containing at least  $n + 1!$  elements and from there deduce that all infinite Boolean algebras satisfy convexity. It is then possible to extend the proof to the case that includes basic, product, and arrow types and deduce the following result:

Every set-theoretic model of closed types in which non-empty types denote infinite sets is a convex set-theoretic model for the polymorphic types.

Therefore, not only do we have a large class of convex models, but also we recover our initial intuition that models with infinite denotations was the way to go.

All that remains to explain is the subtype checking algorithm. We do it in the next section, but before we want to address the possible doubts of a reader about what the denotation of a “finite” type like `Bool` is in such models. In particular, since this denotation contains not only (the denotations of) `true` and `false` but infinitely many other elements, then the reader can rightly wonder what these other elements are and whether they carry any intuitive meaning. In order to explain this point, let us first reexamine what convexity does for infinite types. Convexity is a condition that makes the interpretation of subtyping in some sense independent from the particular syntax of types. Imagine that the syntax of types includes just one basic type: `Int`. Then `Int` is an indivisible type and therefore there exist non-convex models in which the following relation (which is an instance of equation (5) of Section 2.3) holds.

$$\text{Int} \times \alpha \leq (\text{Int} \times \bar{\text{Int}}) \vee (\alpha \times \text{Int}) \quad (9)$$

(*eg*, a model where `Int` is interpreted by a singleton set: in a non-convex model nothing prevents such an interpretation). Now imagine to add the type `Odd`, subtype of `Int`, to the type system: then (9) no longer holds (precisely, the interpretation at issue no longer is a model) since the substitution of  $\alpha$  by `Odd` disproves it. Should the presence of `Odd` change the containment relation between `Int` and the other types? Semantically this should not happen. A relation as (9) should have the same meaning independently from whether `Odd` is included in the syntax of types or not. In other terms we want the addition of `Odd` to yield a conservative extension of the subtyping relation. Therefore, all models in which (9) is valid must be discarded. Convexity does it.

The point is that convexity pushes this idea to all types, so that their interpretation is independent from the possible syntactic subtypes they may have. It is as if the interpretation of subtyping assumed that every type has at least one (actually, infinitely many) stricter non empty subtype(s). So what could the denotation of type

Bool be in such a model, then? A possible choice is to interpret Bool into a set containing labeled versions of true and false, where labels are drawn from an infinite set of labels (a similar interpretation was first introduced by Gesbert *et al.* [14]: see Section 2.8 on related work). Here the singleton type {true} is interpreted as an infinite set containing differently labeled versions of the single element true. Does this labeling carry any intuitive meaning? One can think of it as representing name subtyping: these labels are the names of subtypes of the singleton type {true} for which the subtyping relation is defined by name subtyping. As we do not want the subtyping relation for Int to change (non conservatively) when adding to the system the type Odd, so for the same reason we do not want the subtyping relation for singleton types to change when adding by name subtyping new subtypes, even when these subtypes are subtypes of a singleton type. So convexity makes the subtyping relation insensitive to possible extensions by name subtyping.

## 2.6 Subtyping algorithm

The subtyping algorithm for the relation induced by convex models can be decomposed in 6 elementary steps. Let us explain the intuition underlying each of them: all missing details can be found in Section 3.

First of all, we already said that deciding  $t_1 \leq t_2$ —ie, whether for all  $\eta$ ,  $\llbracket t_1 \rrbracket \eta \subseteq \llbracket t_2 \rrbracket \eta$ —is equivalent to decide the emptiness of the type  $t_1 \wedge \neg t_2$ —ie, whether for all  $\eta$ ,  $\llbracket t_1 \wedge \neg t_2 \rrbracket \eta = \emptyset$ —. So the first step of the algorithm is to transform the problem  $t_1 \leq t_2$  into the problem  $t_1 \wedge \neg t_2 \leq \emptyset$ :

**Step 1:** transform the subtyping problem into an emptiness decision problem.

Our types are just a propositional logic whose atoms are type variables,  $\emptyset$ ,  $\mathbb{1}$ , basic, product, and arrow types. We use  $a$  to range over atoms and, following the logic nomenclature, call *literal*, ranged over by  $\ell$ , an atom or its negation:

$$a ::= b \mid t \times t \mid t \rightarrow t \mid \emptyset \mid \mathbb{1} \mid \alpha \quad \ell ::= a \mid \neg a$$

By using the laws of propositional logic we can transform every type into a disjunctive normal form, that is, into a union of intersections of literals:

$$\bigvee_{i \in I} \bigwedge_{j \in J} \ell_{ij}$$

Since the interpretation function preserves the set-theoretic semantics of type connectives, then every type is empty if and only if its disjunctive normal form is empty. So the second step of our algorithm consists of transforming the type  $t_1 \wedge \neg t_2$  whose emptiness was to be checked, into a disjunctive normal form:

**Step 2:** put the type whose emptiness is to be decided in a disjunctive normal form.

Next, we have to decide when a normal form, that is, a union of intersections, is empty. A union is empty if and only if every member of the union is empty. Therefore the problem reduces to deciding emptiness of an intersection of literals:  $\bigwedge_{i \in I} \ell_i$ . Intersections of literals can be straightforwardly simplified. Every occurrence of the literal  $\mathbb{1}$  can be erased since it does not change the result of the intersection. If either any of the literals is  $\emptyset$  or two literals are a variable and its negation, then we do not have to perform further checks since the intersection is surely empty. An intersection can be simplified also when two literals with different constructors occur in it: if in the intersections there are two atoms of different constructors, say,  $t_1 \times t_2$  and  $t_1 \rightarrow t_2$ , then their intersection is empty and so is the whole intersection; if one of the two atoms is negated, say,  $t_1 \times t_2$  and  $\neg(t_1 \rightarrow t_2)$ , then it can be eliminated since it contains the one that is not negated; if both atoms are negated, then the intersection

can also be simplified (with some more work: cf. the formal development in Section 3). Therefore the third step of the algorithm is to perform these simplifications so that the problem is reduced to deciding emptiness of intersections that are formed by literals that are (possible negations of) either type variables or atoms all of the same constructor (all basic, all product, or all arrow types):

**Step 3:** simplify mixed intersections.

At this stage we have to decide emptiness of intersections of the form  $\bigwedge_{i \in I} a_i \wedge \bigwedge_{j \in J} \neg a'_j \wedge \bigwedge_{h \in H} \alpha_h \wedge \bigwedge_{k \in K} \neg \beta_k$  where all the  $a_i$ 's and  $a'_j$ 's are atoms with the same constructor, and where  $\{\alpha_h\}_{h \in H}$  and  $\{\beta_k\}_{k \in K}$  are disjoint sets of type variables: we just reordered literals so that negated variables and the other negated atoms are grouped together. In this step we want to get rid of the rightmost group in the intersection, that is, the one with negated type variables. In other terms, we want to reduce our problem to deciding the emptiness of an intersections as the above, but where all top-level occurrences of type variables are positive. This is quite easy, and stems from the observation that if a type with a type variable  $\alpha$  is empty for every possible assignment of  $\alpha$ , then it will be empty also if one replaces  $\neg \alpha$  for  $\alpha$  in it: exactly the same checks will be performed since the denotation of the first type for  $\alpha \mapsto S \subseteq \mathcal{D}$  will be equal to the denotation of the second type for  $\alpha \mapsto \bar{S} \subseteq \mathcal{D}$ . That is to say,  $\forall \eta. \llbracket t \rrbracket \eta = \emptyset$  if and only if  $\forall \eta. \llbracket t\{\neg \alpha/\alpha\} \rrbracket \eta = \emptyset$  (where  $t\{\neg \alpha/\alpha\}$  denotes the substitution in  $t$  of  $\neg \alpha$  for  $\alpha$ ). So all the negations of the group of toplevel negated variables can be eliminated by substituting  $\neg \beta_k$  for  $\beta_k$  in the  $a_i$ 's and  $a'_j$ 's:

**Step 4:** eliminate toplevel negative variables.

Next comes what probably is the trickiest step of the algorithm. We have to prove emptiness of intersections of atoms  $a_i$  and negated atoms  $a'_j$  all on the same constructors and of positive variables  $\alpha_k$ . To lighten the presentation let us consider just the case in which atoms are all product types (the case for arrow types is similar though trickier, while the case for basic types is trivial since it reduces to the case for basic types without variables). By using De Morgan's laws we can move negated atoms on the right hand-side of the relation so that we have to check the following containment

$$\bigwedge_{t_1 \times t_2 \in P} t_1 \times t_2 \wedge \bigwedge_{h \in H} \alpha_h \leq \bigvee_{t'_1 \times t'_2 \in N} t'_1 \times t'_2 \quad (10)$$

where  $P$  and  $N$  respectively denote the sets of positive and negative atoms. Our goal is to eliminate all top-level occurrences of variables (the  $\alpha_h$ 's) so that the problem is reduced to checking emptiness of product literals. To that end observe that each  $\alpha_h$  is intersected with other products. Therefore whatever the interpretation of  $\alpha_h$  is, the only part of its denotation that matters is the one that intersects  $\mathcal{D}^2$ . Ergo, it is useless, at least at top-level, to check all possible assignments for  $\alpha_h$ , since those contained in  $\mathcal{D}^2$  will suffice. These can be checked by replacing  $\gamma_h^1 \times \gamma_h^2$  for  $\alpha_h$ , where  $\gamma_h^1, \gamma_h^2$  are fresh type variables. Of course the above reasoning holds for the top-level variables, but nothing tells us that the non top-level occurrences of  $\alpha_h$  will intersect any product. So replacing them with just  $\gamma_h^1 \times \gamma_h^2$  would yield a sound but incomplete check. We rather replace every non toplevel occurrence of  $\alpha_h$  by  $(\gamma_h^1 \times \gamma_h^2) \vee \alpha_h$ . This still is a sound substitution since if (10) holds, then it must also hold for the case where  $(\gamma_h^1 \times \gamma_h^2) \vee \alpha_h$  is substituted for  $\alpha_h$  (with the  $\vee \alpha_h$  part useless for toplevel occurrences). Rather surprisingly, at least at first sight, this substitution is also complete, that is (10) holds if and only if the following holds:

$$\bigwedge_{t_1 \times t_2 \in P} t_1 \theta \times t_2 \theta \wedge \bigwedge_{h \in H} \gamma_h^1 \times \gamma_h^2 \leq \bigvee_{t'_1 \times t'_2 \in N} t'_1 \theta \times t'_2 \theta$$

where  $\theta$  is the substitution  $\{(\gamma_h^1 \times \gamma_h^2) \vee \alpha_h / \alpha_h\}_{h \in H}$ .<sup>5</sup> As an aside, we signal that this transformation holds only because  $\alpha_h$ 's are positive: the application of *Step 4* is thus a necessary precondition to the application of this one. We thus succeeded to eliminate all toplevel occurrences of type variables and, thus, we reduced the initial problem to the problem of deciding emptiness of intersections in which all literals are products or negations of products (and similarly for arrows):

*Step 5: eliminate toplevel variables.*

The final step of our algorithm must decompose the type constructors occurring at toplevel in order to recurse or stop. To that end it will use set-theoretic properties to deconstruct atom types and, above all, the convexity property to decompose the emptiness problem into a set of emptiness subproblems (this is where convexity plays an irreplaceable role: without convexity the definition of an algorithm seems to be out of our reach). Let us continue with our example with products. At this stage all it remains to solve is to decide a containment of the following form (we included the products of fresh variables into  $P$ ):

$$\bigwedge_{t_1 \times t_2 \in P} t_1 \times t_2 \leq \bigvee_{t'_1 \times t'_2 \in N} t'_1 \times t'_2 \quad (11)$$

Using the set-theoretic properties of the interpretation function and our definition of subtyping, we can prove (see Lemma 6.4 in [10] for details) that (11) holds if and only if for all  $N' \subseteq N$ ,

$$\forall \eta. \left( \left[ \bigwedge_{t_1 \times t_2 \in P} t_1 \wedge \bigwedge_{t'_1 \times t'_2 \in N'} \neg t'_1 \right] \eta = \emptyset \text{ or } \left[ \bigwedge_{t_1 \times t_2 \in P} t_2 \wedge \bigwedge_{t'_1 \times t'_2 \in N \setminus N'} \neg t'_2 \right] \eta = \emptyset \right)$$

We can now apply the convexity property and distribute the quantification on  $\eta$  on each subformula of the or. This is equivalent to state that we have to check the emptiness for each type that occurs as argument of the interpretation function. Playing a little more with De Morgan's laws and applying the definition of subtyping we can thus prove that (11) holds if and only if

$$\forall N' \subseteq N. \left( \bigwedge_{t_1 \times t_2 \in P} t_1 \leq \bigvee_{t'_1 \times t'_2 \in N'} t'_1 \right) \text{ or } \left( \bigwedge_{t_1 \times t_2 \in P} t_2 \leq \bigvee_{t'_1 \times t'_2 \in N \setminus N'} t'_2 \right)$$

To understand the rationale of this transformation the reader can consider the case in which both  $P$  and  $N$  contain just one atom, namely, the case for  $t_1 \times t_2 \leq t'_1 \times t'_2$ . There are just two cases to check ( $N' = \emptyset$  and  $N' = N$ ) and it is not difficult to see that the condition above becomes:  $(t_1 \leq 0)$  or  $(t_2 \leq 0)$  or  $(t_1 \leq t'_1 \text{ and } t_2 \leq t'_2)$ , as expected.

The important point however is that we were able to express the problem of (11) in terms of subproblems that rest on strict subterms (there is a similar decomposition rule for arrow types). Remember that our types are possibly infinite trees since they were *coinductively* generated by the grammar in §2.1. We do not consider every possible coinductively generated tree, but only those that are regular (*ie*, that have a finite number of distinct subtrees) and in which every infinite branch contains infinitely many occurrences of type constructors (*ie*, products and arrows). The last condition rules out meaningless terms (such as  $t = \neg t$ ) as well as infinite unions and intersections. It also provides a well-founded order that allows us to use recursion. Therefore, we memoize the relation in (11) and recursively call the algorithm from *Step 1* on the subterms we obtained from decomposing the toplevel constructors:

<sup>5</sup>Note that the result of this substitution is equivalent to using the substitution  $\{(\gamma_h^1 \times \gamma_h^2) \vee \gamma_h^3 / \alpha_h\}_{h \in H}$  where  $\gamma_h^3$  is also a fresh variable: we just spare a new variable by reusing  $\alpha_h$  which would be no longer used (actually this artifice makes proofs much easier).

*Step 6: eliminate toplevel constructors, memoize, and recurse.*

The algorithm is sound and complete with respect to the subtyping relation defined by (7) and terminates on all types (which implies the decidability of the subtyping relation).

## 2.7 Examples

The purpose of this subsection is twofold: first, we want to give some examples to convey the idea that the subtyping relation is intuitive; second we present some cases that justify the subtler and more technical aspects of the subtyping algorithm we exposed in the previous subsection. All the examples below can be tested in our prototype subtype-checker.

In what follows we will use  $x, y, z$  to range over recursion variables and the notation  $\mu x.t$  to denote recursive types. This should suffice to avoid confusion with free type variables that are ranged over by  $\alpha, \beta$ , and  $\gamma$ .

As a first example we show how to use type variables to internalize meta-properties. For instance, for all ground types  $t_1, t_2$ , and  $t_3$  the relation  $(t_1 \rightarrow t_3) \wedge (t_2 \rightarrow t_3) \leq (t_1 \vee t_2) \rightarrow t_3$  and its converse hold. This meta-theoretic property can be expressed in our type system since the following relation holds:

$$(\alpha \rightarrow \gamma) \wedge (\beta \rightarrow \gamma) \sim (\alpha \vee \beta) \rightarrow \gamma$$

(where  $\sim$  denotes that both  $\leq$  and  $\geq$  hold). Of course we can apply this generalization to any relation that holds for generic types. For instance, we can prove common distributive laws such as

$$((\alpha \vee \beta) \times \gamma) \sim (\alpha \times \gamma) \vee (\beta \times \gamma) \quad (12)$$

and combine it with the previous relation and the covariance of arrow on codomains to deduce

$$(\alpha \times \gamma \rightarrow \delta_1) \wedge (\beta \times \gamma \rightarrow \delta_2) \leq ((\alpha \vee \beta) \times \gamma) \rightarrow \delta_1 \vee \delta_2$$

Similarly we can prove that the set of lists whose elements have type  $\alpha$ , that is,  $\alpha \text{ list} = \mu x.(\alpha \times x) \vee \text{nil}$ , contains both the  $\alpha$ -lists with an even number of elements

$$\mu x.(\alpha \times (\alpha \times x)) \vee \text{nil} \leq \mu x.(\alpha \times x) \vee \text{nil}$$

(where  $\text{nil}$  denotes the singleton type containing just the value  $\text{nil}$ ) and the  $\alpha$ -lists with an odd number of elements

$$\mu x.(\alpha \times (\alpha \times x)) \vee (\alpha \times \text{nil}) \leq \mu x.(\alpha \times x) \vee \text{nil}$$

and it is itself contained in the union of the two, that is:

$$\alpha \text{ list} \sim (\mu x.(\alpha \times (\alpha \times x)) \vee \text{nil}) \vee (\mu x.(\alpha \times (\alpha \times x)) \vee (\alpha \times \text{nil}))$$

We said that the intuition for subtyping type variables is to consider them as basic types. But type variables are *not* basic types. As an example, if  $t$  is a non-empty type, then we have that:

$$\alpha \wedge (\alpha \times t) \not\leq t_1 \rightarrow t_2$$

which implies that  $\alpha \wedge (\alpha \times t)$  is not empty. This is correct because if for instance we substitute the type  $t \vee (t \times t)$  for  $\alpha$ , then (by the distributivity law stated in (12)) the intersection is equal to  $(t \times t)$ , which is non-empty. However, note that if  $\alpha$  were a basic type, then the intersection  $\alpha \wedge (\alpha \times t)$  would be empty. Furthermore, since the following relation holds

$$\alpha \wedge (\alpha \times t) \leq \alpha$$

then, this last containment is an example of non-trivial containment (in the sense that the left hand-side is not empty) involving type variables. For an example of non-trivial containment involving arrows the reader can check

$$\mathbb{1} \rightarrow \mathbb{0} \leq \alpha \rightarrow \beta \leq \mathbb{0} \rightarrow \mathbb{1}$$

which states that  $\mathbb{1} \rightarrow \mathbb{0}$ , the set of all functions that diverge on all arguments, is contained in all arrow types  $\alpha \rightarrow \beta$  (whatever types

$\alpha$  and  $\beta$  are) and that the latter are contained in  $\mathbb{0} \rightarrow \mathbb{1}$ , which is the set of all function values.

Type connectives implement classic proposition logic. If we use  $\alpha \Rightarrow \beta$  to denote  $\neg\alpha \vee \beta$ , that is logical implication, then the following subtyping relation is a proof of Pierce’s law:

$$\mathbb{1} \leq ((\alpha \Rightarrow \beta) \Rightarrow \alpha) \Rightarrow \alpha$$

since being a supertype of  $\mathbb{1}$  logically corresponds to being equivalent to true (note that arrow types do not represent logical implication; for instance,  $\mathbb{0} \rightarrow \mathbb{1}$  is not empty: it contains all function values). Similarly, the system captures the fundamental property that for all non-empty sets  $\beta$  the set  $(\beta \wedge \alpha) \vee (\beta \wedge \neg\alpha)$  is never empty:

$$(\beta \wedge \alpha) \vee (\beta \wedge \neg\alpha) \sim \beta$$

from which we can derive

$$\mathbb{1} \leq (((\beta \wedge \alpha) \vee (\beta \wedge \neg\alpha)) \Rightarrow \mathbb{0}) \Rightarrow (\beta \Rightarrow \mathbb{0})$$

This last relation can be read as follows: if  $(\beta \wedge \alpha) \vee (\beta \wedge \neg\alpha)$  is empty, then  $\beta$  is empty.

But the property above will never show a stuttering validity since the algorithm returns false when asked to prove

$$\text{nil} \times \alpha \leq (\text{nil} \times \neg\text{nil}) \vee (\alpha \times \text{nil})$$

even for a singleton type as nil.

The subtyping relation has some simple form of introspection since  $t_1 \leq t_2$  if and only if  $\mathbb{1} \leq t_1 \Rightarrow t_2$  (ie, by negating both types and reversing the subtyping relation,  $t_1 \wedge \neg t_2 \leq \mathbb{0}$ ). However, the introspection capability is very limited insofar as it is possible to state interesting properties only when atoms are type variables: although we can characterize the subtyping relation  $\leq$ , we have no idea about how to characterize its negation  $\not\leq$ .<sup>6</sup>

The necessity for the tricky substitution  $\alpha \mapsto (\gamma_1 \times \gamma_2) \vee \alpha$  performed at *Step 5* of the algorithm can be understood by considering the following example where  $t$  is any non-empty type:

$$(\alpha \times t) \wedge \alpha \leq ((\mathbb{1} \times \mathbb{1}) \times t).$$

If in order to check the relation above we substituted just  $\gamma_1 \times \gamma_2$  for  $\alpha$ , then this would yield a positive result, which is wrong: if we replace  $\alpha$  by  $b \vee (b \times t)$ , where  $b$  is any basic type, then the intersection on the left becomes  $(b \times t)$  and  $b$  is neither contained in  $\mathbb{1} \times \mathbb{1}$  nor empty. Our algorithm correctly disproves the containment, since it checks also the substitution of  $(\gamma_1 \times \gamma_2) \vee \alpha$  for the first occurrence of  $\alpha$ , which captures the above counterexample.

Finally, the system also proves subtler relations whose meaning is not clear at first sight, such as:

$$\alpha_1 \rightarrow \beta_1 \leq ((\alpha_1 \wedge \alpha_2) \rightarrow (\beta_1 \wedge \beta_2)) \vee \neg(\alpha_2 \rightarrow (\beta_2 \wedge \neg\beta_1)) \quad (13)$$

In order to prove it, the subtyping algorithm first moves the occurrence of  $\alpha_2 \rightarrow (\beta_2 \wedge \neg\beta_1)$  from the right of the subtyping relation to its left:  $(\alpha_1 \rightarrow \beta_1) \wedge (\alpha_2 \rightarrow (\beta_2 \wedge \neg\beta_1)) \leq ((\alpha_1 \wedge \alpha_2) \rightarrow (\beta_1 \wedge \beta_2))$ ; then following the decomposition rules for arrows the algorithm checks the four following cases (Step 6 of the algorithm), which hold straightforwardly:

$$\left\{ \begin{array}{l} \alpha_1 \wedge \alpha_2 \leq \mathbb{0} \text{ or } \beta_1 \wedge (\beta_2 \wedge \neg\beta_1) \leq \beta_1 \wedge \beta_2 \\ \alpha_1 \wedge \alpha_2 \leq \alpha_1 \text{ or } (\beta_2 \wedge \neg\beta_1) \leq \beta_1 \wedge \beta_2 \\ \alpha_1 \wedge \alpha_2 \leq \alpha_2 \text{ or } \beta_1 \leq \beta_1 \wedge \beta_2 \\ \alpha_1 \wedge \alpha_2 \leq \alpha_1 \vee \alpha_2 \end{array} \right.$$

<sup>6</sup> For instance, it would be nice to prove something like:

$$(\neg\beta_1 \vee ((\beta_1 \Rightarrow \alpha_1) \wedge (\alpha_2 \Rightarrow \beta_2))) \sim (\alpha_1 \rightarrow \alpha_2 \Rightarrow \beta_1 \rightarrow \beta_2)$$

since it seems to provide a complete characterization of the subtyping relation between two arrow types. Unfortunately the equivalence is false since  $\beta_1 \not\leq \alpha_1$  does not imply  $\beta_1 \wedge \neg\alpha_1 \geq \mathbb{1}$  but just  $\beta_1 \wedge \neg\alpha_1 \leq \mathbb{0}$ . This property can be stated only at meta level, that is:  $\alpha_1 \rightarrow \alpha_2 \leq \beta_1 \rightarrow \beta_2$  if and only if  $(\beta_1 \leq \mathbb{0} \text{ or } (\beta_1 \leq \alpha_1 \text{ and } \alpha_2 \leq \beta_2))$ .

Notice that relation (13) is quite subtle insofar as neither  $\alpha_1 \rightarrow \beta_1 \leq (\alpha_1 \wedge \alpha_2) \rightarrow (\beta_1 \wedge \beta_2)$  nor  $\alpha_1 \rightarrow \beta_1 \leq \neg(\alpha_2 \rightarrow (\beta_2 \wedge \neg\beta_1))$  hold: the type on left hand-side of (13) is contained in the union of the two types on the right hand-side of (13) without being completely contained in either of them.

## 2.8 Related work

This work extends the work on semantic subtyping [10], as such the two works share the same approach and common developments. Since convex models of our theory can be derived from the models of [10], then several techniques we used in our subtyping algorithm (in particular the decomposition of toplevel type constructors) are directly issued from the research in [10]. This work starts precisely from where [10] stopped, that is the monomorphic case, and adds prenex parametric polymorphism to it.

The most advanced work on polymorphism for XML types, thus far, is the Hosoya, Frisch, and Castagna’s approach described in [15], whose extended abstract was first presented at POPL ’05. Together with [10], the paper by Hosoya, Frisch, and Castagna constitutes the starting point of this work. A starting point that, so far, was rather considered to establish a (negative) final point. As a matter of fact, although the polymorphic system in [15] is the one used to define the polymorphic extension of XDuce [16] (incorporated from version 0.5.0 of the language), the three authors of [15] agree that the main interest of their work does not reside in its type system, but rather in the negative results that motivate it. In particular, the pivotal example of our work, equation (5), was first presented in [15], and used there to corroborate the idea that a purely semantic approach for polymorphism of regular tree types was an hopeless quest. At that time, this seemed so more hopeless that the equation (5) did not involve arrow types: a semantically defined polymorphic subtyping looked out of reach even in the restrictive setting of Hosoya and Pierce seminal work [16], which did not account for higher-order functions. This is why [15] falls back on a syntactic approach that, even if it retains some flavors of semantic subtyping, it cannot be extended to higher-order functions (a lack that nevertheless fits XDuce). Our works shows that the negative results of [15] were not so insurmountable as it had been thought.

Hitherto, the only work that blends polymorphic regular types and arrow types is Jérôme Vouillon’s work that was presented at POPL ’06 [20]. His approach, however, is very different from ours insofar as it is intrinsically syntactic. Vouillon starts from a particular language (actually, a peculiar pattern algebra) and coinductively builds up on it the subtyping relation by a set of inference rules. The type algebra includes only the union connective (negation and intersection are missing) and a semantic interpretation of subtyping is given *a posteriori* by showing that a pattern (types are special cases of patterns) can be considered as the set of values that match the pattern. Nevertheless, this interpretation is still syntactic in nature since it relies on the definition of matching and containment, yielding a system tailored for the peculiar language of the paper. This allows Vouillon to state impressive and elegant results such as the translation of the calculus into a non-explicitly-typed one, or the interpretation of open types containment as in our equation (4) (according to Vouillon this last result is made possible in his system by the absence of intersection types, although the critical example in (5) does not involve any intersection). But the price to pay is a system that lacks naturalness (eg, the wild-card pattern has different meanings according to whether it occurs in the right or in left type of a subtyping relation) and, even more, it lacks the generality of our approach (we did not state our subtype system for any specific language while Vouillon’s system is inherently tied to a particular language whose semantics it completely relies on). The semantics

of Vouillon’s patterns is so different from ours that typing Vouillon’s language with our types seems quite difficult.

Other works less related to ours are those in which XML and polymorphism are loosely coupled. This is the case of OCaml-Duce [9] where ML-polymorphism and XML types and patterns are merged together without mixing: the main limitation of this approach is that it does not allow parametric polymorphism for XML types, which is the whole point of our (and Vouillon’s) work(s). A similar remark can be done for Xtatic [11] that merges C# name subtyping with the XDuce set-theoretic subtyping and for XHaskell [17] whose main focus is to implement XML subtyping using Haskell’s type-classes. A more thorough comparison of these approaches can be found in [9, 15].

Polymorphism can be attained by adopting the so-called data-binding approach which consists in encoding XML types and values into the structures of an existing polymorphic programming language. This is the approach followed by HaXML [23]. While the polymorphism is inherited from the target language, the rigid encoding of XML data into fixed structures loses all flexibility of the XML type equivalences so as, for instance,  $(t \times s_1) \vee (t \times s_2)$  and  $(t \times s_1 \vee s_2)$  are different (and even unrelated) types.

Our work already has a follow-up. In a paper included in these same proceedings [14] Gesbert, Genevès, and Layaida use the framework we define here to give a different decision procedure for our subtyping relation. More precisely, they take a specific model for the monomorphic type system (ie, the model defined by Frisch *et al.* [10] and used by the language CDuce), they encode the subtyping relation induced by this model into a tree logic, and use a satisfiability solver to efficiently decide it. Next, they extend the type system with type variables and they obtain a convex model by interpreting non-empty types as infinite sets using a labeling technique similar to the one we outlined at the end of Section 2.5: they label values by (finite sets of) type variables and every non empty ground type is, thus, interpreted as an infinite set containing the infinitely many labelings of its values. Again the satisfiability solver provides a decision procedure for the subtyping relation. Their technique is interesting in several respects. First it provides a very elegant solution to the problem of deciding our subtyping relation, solution that is completely different from the one given here. Second, their technique shows that the decision problem is EXPTIME, (while here we only prove the decidability of the problem by showing the termination of our algorithm). Finally, their logical encoding paves the way to extending types (and subtyping) with more expressive logical constraints representable by their tree logic. In contrast, our algorithm is interesting for quite different reasons: first, it is defined for generic interpretations rather than for a fixed model; second, it shows how convexity is used in practice (see in particular **Step 6** of the algorithm); and, finally, our algorithm is a straightforward modification of the algorithm used in CDuce and, as such, can benefit of the technology and optimizations used there.<sup>7</sup> We expect the integration of this subtyping relation in the CDuce to be available in the near future.

Finally, we signal the work on polymorphic iterators for XML presented in [4]. It introduces a very simple strongly normalizing calculus fashioned to define tree iterators. These iterators are lightly checked at the moment of their definition: the compiler does not complain unless they are irremediably flawed. This optimistic typing, combined with the relatively limited expressive power of the calculus, makes it possible to type iterator applications in a very precise way (essentially, by performing an abstract execution of the iterator on the types) yielding a kind of polymorphism that

<sup>7</sup> Alain Frisch’s PhD. thesis [8] describes two algorithms that improve over the simple saturation-based strategy described in Section 3.4. They are used both in CDuce compiler and in the prototype we implemented to check the subtyping relation presented in this work.

is out of reach of parametric or subtype polymorphism (for instance it can precisely type the reverse function applied to *heterogeneous* lists and thus deduce that the application of reverse to a list of type, say,  $[ \text{Int Bool}^* \text{Char}^+ ]$  yields a result of type  $[ \text{Char}^+ \text{Bool}^* \text{Int} ]$ ). As such it is orthogonal to the kind of polymorphism presented here, and both can and should coexist in a same language.

### 3. Formal development

In this section we describe the technical development that supports the results we exposed in the previous section. We (strongly) suggest readers to skip this section on first reading and directly jump to the conclusions in Section 4. For space reasons all proofs are omitted. They can be found in the extended version available from the first author’s web page together with a link to the source code of our subtype checker prototype. The prose is reduced to the strict necessary, and limited to explain points that were not dealt with in the previous section.

#### 3.1 Types

**Definition 3.1. (Types)** Consider a countable set of type variables  $\mathcal{V}$  ranged over by  $\alpha$  and a finite set of basic (or constant) types ranged over by  $b$ . A type is a regular tree coinductively produced by the following grammar

$$t ::= \alpha \mid b \mid t \times t \mid t \rightarrow t \mid t \vee t \mid \neg t \mid \mathbb{0}$$

and in which every infinite branch contains infinitely many occurrences of atoms (ie, either a type variable or the immediate application of a type constructor: basic, product, arrow)

We write  $t_1 \setminus t_2$ ,  $t_1 \wedge t_2$ , and  $\mathbb{1}$  respectively as abbreviation for  $t_1 \wedge \neg t_2$ ,  $\neg(\neg t_1 \vee \neg t_2)$ , and  $\neg \mathbb{0}$ .

The condition on infinite branches bars out ill-formed types such as  $t = t \vee t$  (which does not carry any information about the set denoted by the type) or  $t = \neg t$  (which cannot represent any set). It also ensures that the binary relation  $\triangleright \subseteq \mathcal{T}^2$  defined by  $t_1 \vee t_2 \triangleright t_i$ ,  $\neg t \triangleright t$  is Noetherian (that is, strongly normalizing). This gives an induction principle on  $\mathcal{T}$  that we will use without any further explicit reference to the relation.

**Notation.** Let  $t$  be a type. We use  $\text{var}(t)$  to denote the set of type variables occurring in  $t$  and by  $\text{tlv}(t)$  (toplevel variables) all the variables of  $t$  that have at least one occurrence not under a constructor, that is:  $\text{tlv}(\alpha) = \{\alpha\}$ ,  $\text{tlv}(\neg t) = \text{tlv}(t)$ ,  $\text{tlv}(t \vee s) = \text{tlv}(t) \cup \text{tlv}(s)$ , and  $\text{tlv}(t) = \emptyset$  otherwise. We say that  $t$  is ground or closed if and only if  $\text{var}(t)$  is empty.

#### 3.2 Subtyping

**Definition 3.2 (Set-Theoretic Interpretation).** A set-theoretic interpretation of  $\mathcal{T}$  is given by a set  $\mathcal{D}$  and a function  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  such that, for all  $t_1, t_2, t \in \mathcal{T}$ ,  $\alpha \in \mathcal{V}$  and  $\eta \in \mathcal{P}(\mathcal{D})^\vee$ :  $\llbracket t_1 \vee t_2 \rrbracket \eta = \llbracket t_1 \rrbracket \eta \cup \llbracket t_2 \rrbracket \eta$ ,  $\llbracket \neg t \rrbracket \eta = \mathcal{D} \setminus \llbracket t \rrbracket \eta$ ,  $\llbracket \alpha \rrbracket \eta = \eta(\alpha)$ , and  $\llbracket \mathbb{0} \rrbracket \eta = \emptyset$ .

**Definition 3.3. (Subtyping Relation)** Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a set-theoretic interpretation. We define the subtyping relation  $\leq_{\llbracket \_ \rrbracket} \subseteq \mathcal{T}^2$  as follows:

$$t \leq_{\llbracket \_ \rrbracket} s \Leftrightarrow \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \llbracket t \rrbracket \eta \subseteq \llbracket s \rrbracket \eta$$

We write  $t \leq s$  when the interpretation  $\llbracket \_ \rrbracket$  is clear from the context.

For each basic type  $b$ , we assume there is a fixed set of constants  $\mathbb{B}(b) \subseteq \mathcal{C}$  whose elements are called constants of type  $b$ . For two basic types  $b_1, b_2$ , the sets  $\mathbb{B}(b_i)$  can have a non-empty intersection.

If, as suggested in Section 2, we interpret extensionally an arrow  $t_1 \rightarrow t_2$  as  $\mathcal{P}(\llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket)$  (precisely as  $\mathcal{P}(\mathcal{D}^2 \setminus (\llbracket t_1 \rrbracket \times (\mathcal{D} \setminus \llbracket t_2 \rrbracket)))$ ),



then every function type is a subtype of  $\mathbb{1} \rightarrow \mathbb{1}$ . We do not want such a property to hold because, otherwise, we could subsume every function to a function that accepts every value and, therefore, every application of a well-typed function to a well-typed argument would be well-typed, independently from the types of the function and of the argument. For example, if, say,  $\text{succ} : \text{Int} \rightarrow \text{Int}$ , then we could deduce  $\text{succ} : \mathbb{1} \rightarrow \mathbb{1}$  and then  $\text{succ}(\text{true})$  would have type  $\mathbb{1}$ . To avoid this problem we introduce an explicit type error  $\Omega$  and use it to define function spaces:

**Definition 3.4.** If  $D$  is a set and  $X, Y$  are subsets of  $D$ , we write  $D_\Omega$  for  $D + \{\Omega\}$  and define  $X \rightarrow Y$  as:

$$X \rightarrow Y \stackrel{\text{def}}{=} \{f \subseteq D \times D_\Omega \mid \forall (d, d') \in f. d \in X \Rightarrow d' \in Y\}$$

This is used in the definition of the extensional interpretation:

**Definition 3.5 (Extensional Interpretation).** Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a set-theoretic interpretation. Its associated extensional interpretation is the unique function  $\mathbb{E}(\_ ) : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathbb{E}\mathcal{D})$  where  $\mathbb{E}\mathcal{D} = \mathcal{D} + \mathcal{D}^2 + \mathcal{P}(\mathcal{D} \times \mathcal{D}_\Omega)$ , defined as follows:

$$\begin{aligned} \mathbb{E}(\alpha)\eta &= \eta(\alpha) \subseteq \mathcal{D} \\ \mathbb{E}(b)\eta &= \mathbb{B}(b) \subseteq \mathcal{D} \\ \mathbb{E}(t_1 \times t_2)\eta &= \llbracket t_1 \rrbracket \eta \times \llbracket t_2 \rrbracket \eta \subseteq \mathcal{D}^2 \\ \mathbb{E}(t_1 \rightarrow t_2)\eta &= \llbracket t_1 \rrbracket \eta \rightarrow \llbracket t_2 \rrbracket \eta \subseteq \mathcal{P}(\mathcal{D} \times \mathcal{D}_\Omega) \\ \mathbb{E}(0)\eta &= \emptyset \\ \mathbb{E}(t_1 \vee t_2)\eta &= \mathbb{E}(t_1)\eta \cup \mathbb{E}(t_2)\eta \\ \mathbb{E}(\neg t)\eta &= \mathbb{E}\mathcal{D} \setminus \mathbb{E}(t)\eta \end{aligned}$$

The definition of set-theoretic model is then as expected:

**Definition 3.6 (Foundation, Convexity, and Models).** Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a set-theoretic interpretation. It is

1. convex if it satisfies equation (8) for all finite choices of types  $t_1, \dots, t_n$ ;
2. structural if  $\mathcal{D}^2 \subseteq \mathcal{D}$ ,  $\llbracket t_1 \times t_2 \rrbracket \eta = \llbracket t_1 \rrbracket \eta \times \llbracket t_2 \rrbracket \eta$  and the relation on  $\mathcal{D}$  induced by  $(d_1, d_2) \blacktriangleright d_i$  is Noetherian;
3. a model if it induces the same subtyping relation as its associated extensional interpretation, that is:  $\forall t \in \mathcal{T}. \forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \llbracket t \rrbracket \eta = \emptyset \iff \mathbb{E}(t)\eta = \emptyset$ . A model is convex if its set-theoretic interpretation is convex; it is well-founded if it induces the same subtyping relation as a structural set-theoretic interpretation.

From now on we consider only well-founded convex models. We already explained the necessity of the notion of convexity we introduced in §2.5. The notion of well-founded model was introduced in §4.3 of [10]. Intuitively, well-founded models are models that contain only values that are finite (eg, in a well-founded model the type  $\mu x.(x \times x)$ —ie, the type that “should” contain all and only infinite binary trees—is empty). This fits the practical motivations of this paper, since XML documents—ie, values—are finite trees.

### 3.3 Properties of the subtyping relation

We write  $\mathcal{A}_{\text{fun}}$  for atoms of the form  $t_1 \rightarrow t_2$ ,  $\mathcal{A}_{\text{prod}}$  for atoms of the form  $t_1 \times t_2$ ,  $\mathcal{A}_{\text{basic}}$  for basic types, and  $\mathcal{A}$  for  $\mathcal{A}_{\text{fun}} \cup \mathcal{A}_{\text{prod}} \cup \mathcal{A}_{\text{basic}}$ . Therefore  $\mathcal{V} \cup \mathcal{A} \cup \{\emptyset, \mathbb{1}\}$  denotes the set of all atoms. Henceforth, we will disregard the atoms  $\emptyset$  and  $\mathbb{1}$  since they can be straightforwardly eliminated during the algorithmic treatment of subtyping.

**Definition 3.7. (Normal Form)** A (disjunctive) normal form  $\tau$  is a finite set of pairs of finite sets of atoms, that is, an element of  $\mathcal{P}_f(\mathcal{P}_f(\mathcal{A} \cup \mathcal{V}) \times \mathcal{P}_f(\mathcal{A} \cup \mathcal{V}))$  (where  $\mathcal{P}_f(\_)$  denotes the finite powerset). Moreover, we call an element of  $\mathcal{P}_f(\mathcal{A} \cup \mathcal{V}) \times \mathcal{P}_f(\mathcal{A} \cup \mathcal{V})$

$\mathcal{V}$ ) a single normal form. If  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  is an arbitrary set-theoretic interpretation,  $\tau$  a normal form and  $\eta$  an assignment, we define  $\llbracket \tau \rrbracket \eta$  as:

$$\llbracket \tau \rrbracket \eta = \bigcup_{(P, N) \in \tau} \bigcap_{a \in P} \llbracket a \rrbracket \eta \cap \bigcap_{a \in N} (\mathcal{D} \setminus \llbracket a \rrbracket \eta)$$

(with the convention that an intersection over an empty set is taken to be  $\mathcal{D}$ ).

For every type  $t$  there exists a normal form  $\mathcal{N}(t)$  that has the same interpretation, and vice versa. For instance, consider the type  $t = a_1 \wedge (a_3 \vee \neg a_2)$  where  $a_1, a_2$ , and  $a_3$  are any atoms. Then  $\mathcal{N}(t) = \{(\{a_1, a_3\}, \emptyset), (\{a_1\}, \{a_2\})\}$ . This corresponds to the fact that for every set-theoretic interpretation and semantic assignment  $\mathcal{N}(t)$ ,  $t$ , and  $(a_1 \wedge a_3) \vee (a_1 \wedge \neg a_2)$  have the same denotation.

For these reasons henceforth we will often confound the notions of types and normal forms, and often speak of the “type”  $\tau$ , taking the latter as a canonical representation of all the types in  $\mathcal{N}^{-1}(\tau)$ .

Let  $\llbracket \_ \rrbracket$  be a set-theoretic interpretation. Given a normal form  $\tau$ , we are interested in comparing the assertions  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \mathbb{E}(\tau)\eta = \emptyset$  and  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \llbracket \tau \rrbracket \eta = \emptyset$ . Clearly, the equation  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \mathbb{E}(\tau)\eta = \emptyset$  is equivalent to:

$$\forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \forall (P, N) \in \tau. \bigcap_{a \in P} \mathbb{E}(a)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \quad (14)$$

Let us write  $\mathbb{E}^{\text{basic}}\mathcal{D} = \mathcal{D}$ ,  $\mathbb{E}^{\text{prod}}\mathcal{D} = \mathcal{D}^2$  and  $\mathbb{E}^{\text{fun}}\mathcal{D} = \mathcal{P}(\mathcal{D} \times \mathcal{D}_\Omega)$ . Then we have  $\mathbb{E}\mathcal{D} = \bigcup_{u \in U} \mathbb{E}^u\mathcal{D}$  where  $U = \{\text{basic}, \text{prod}, \text{fun}\}$ . Thus we can rewrite Inequality (14) as:

$$\forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \forall u \in U. \forall (P, N) \in \tau. \bigcap_{a \in P} (\mathbb{E}(a)\eta \cap \mathbb{E}^u\mathcal{D}) \subseteq \bigcup_{a \in N} (\mathbb{E}(a)\eta \cap \mathbb{E}^u\mathcal{D}) \quad (15)$$

For an atom  $a \in \mathcal{A}$ , we have  $\mathbb{E}(a)\eta \cap \mathbb{E}^u\mathcal{D} = \emptyset$  if  $a \notin \mathcal{A}_u$  and  $\mathbb{E}(a)\eta \cap \mathbb{E}^u\mathcal{D} = \mathbb{E}(a)\eta$  if  $a \in \mathcal{A}_u$ . Then we can rewrite Inequality (15) as:

$$\begin{aligned} \forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \forall u \in U. \forall (P, N) \in \tau. (P \subseteq \mathcal{A}_u \cup \mathcal{V}) \Rightarrow \\ \bigcap_{a \in P \cap \mathcal{A}_u} \mathbb{E}(a)\eta \cap \bigcap_{\alpha \in P \cap \mathcal{V}} (\eta(\alpha) \cap \mathbb{E}^u\mathcal{D}) \subseteq \\ \bigcup_{a \in N \cap \mathcal{A}_u} \mathbb{E}(a)\eta \cup \bigcup_{\alpha \in N \cap \mathcal{V}} (\eta(\alpha) \cap \mathbb{E}^u\mathcal{D}) \end{aligned} \quad (16)$$

(where the intersection is taken to be  $\mathbb{E}\mathcal{D}$  when  $P = \emptyset$ ). Furthermore, if the same variable occurs both in  $P$  and in  $N$ , then (16) is trivially satisfied. So we can suppose that  $P \cap N \cap \mathcal{V} = \emptyset$ . This justifies the simplifications made in Step 3 of the subtyping algorithm, that is the simplification of mixed single normal forms.

Step 4, the elimination of negated toplevel variables, is justified by the following lemma:

**Lemma 3.8.** Let  $P, N$  be two finite subsets of atoms and  $\alpha$  an arbitrary type variable, then we have

$$\begin{aligned} \forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \bigcap_{a \in P} \mathbb{E}(a)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \cup \eta(\alpha) \Leftrightarrow \\ \forall \eta \in \mathcal{P}(\mathcal{D})^\vee. \bigcap_{a \in P} \mathbb{E}(\theta_\alpha^{-\beta}(a))\eta \cap \eta(\beta) \subseteq \bigcup_{a \in N} \mathbb{E}(\theta_\alpha^{-\beta}(a))\eta \end{aligned}$$

where  $\beta$  is a fresh variable and  $\theta_\alpha^{-\beta}(a) = a\{\neg\beta/\alpha\}$

Note that Lemma 3.8 only deals with one type variable, but it is trivial to generalize this lemma to multiple type variables (the same holds for Lemmas 3.9 and 3.10).

Using Lemma 3.8, we can rewrite Inequality (16) as:

$$\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \forall u \in U . \forall (P, N) \in \tau . (P \subseteq \mathcal{A}_u \cup \mathcal{V}) \Rightarrow \bigcap_{a \in P \cap \mathcal{A}_u} \mathbb{E}(a)\eta \cap \bigcap_{\alpha \in P \cap \mathcal{V}} (\eta(\alpha) \cap \mathbb{E}^u \mathcal{D}) \subseteq \bigcup_{a \in N \cap \mathcal{A}_u} \mathbb{E}(a)\eta \quad (17)$$

since we can assume  $N \cap \mathcal{V} = \emptyset$ .

Next, we justify *Step 5* of the algorithm, that is the elimination of the toplevel variables. In (17) this corresponds to eliminating the variables in  $P \cap \mathcal{V}$ . When  $u = \text{basic}$  this can be easily done since all variables (which can appear only at top-level) can be straightforwardly removed. Indeed, notice that if  $s$  and  $t$  are closed types then  $s \wedge \alpha \leq t$  if and only if  $s \leq t$ .

The justification of *Step 5* for  $u = \text{prod}$  is given by Lemma 3.9.

**Lemma 3.9.** *Let  $(\mathcal{D}, \llbracket \cdot \rrbracket)$  be a well-founded model,  $P, N$  two finite subsets of  $\mathcal{A}_{\text{prod}}$  and  $\alpha$  an arbitrary type variable.*

$$\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \bigcap_{a \in P} \mathbb{E}(a)\eta \cap (\eta(\alpha) \cap \mathbb{E}^{\text{prod}} \mathcal{D}) \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \Leftrightarrow \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \bigcap_{a \in P} \mathbb{E}(\theta_\alpha^\times(a))\eta \cap \mathbb{E}(\alpha_1 \times \alpha_2)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(\theta_\alpha^\times(a))\eta$$

where  $\theta_\alpha^\times(a) = a\{(\alpha_1 \times \alpha_2) \vee \alpha / \alpha\}$  and  $\alpha_1, \alpha_2$  are fresh variables.

The case for  $u = \text{fun}$  is trickier because  $\mathbb{1} \rightarrow \mathbb{0}$  is contained in every arrow type, and therefore sets of arrows that do not contain it must be explicitly checked. If, analogously to  $u = \text{prod}$ , we used just  $\{(\alpha_1 \rightarrow \alpha_2) \vee \alpha / \alpha\}$ , then the subtypes of  $\mathbb{0} \rightarrow \mathbb{1}$  that do not contain  $\mathbb{1} \rightarrow \mathbb{0}$  would never be checked to  $\alpha$  by the algorithm and, thus, never checked. Therefore,  $\{((\alpha_1 \rightarrow \alpha_2) \setminus (\mathbb{1} \rightarrow \mathbb{0})) \vee \alpha / \alpha\}$  must be checked, as well.

**Lemma 3.10.** *Let  $(\mathcal{D}, \llbracket \cdot \rrbracket)$  be a well-founded model,  $P, N$  two finite subsets of  $\mathcal{A}_{\text{fun}}$  and  $\alpha$  an arbitrary type variable. Then*

$$\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \bigcap_{a \in P} \mathbb{E}(a)\eta \cap (\eta(\alpha) \cap \mathbb{E}^{\text{fun}} \mathcal{D}) \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \Leftrightarrow \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \bigcap_{a \in P} \mathbb{E}(\theta_\alpha^{\rightarrow}(a))\eta \cap \mathbb{E}(\alpha_1 \rightarrow \alpha_2)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(\theta_\alpha^{\rightarrow}(a))\eta$$

and  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \bigcap_{a \in P} \mathbb{E}(\theta_\alpha^{\rightsquigarrow}(a))\eta \cap \mathbb{E}((\alpha_1 \rightarrow \alpha_2) \setminus (\mathbb{1} \rightarrow \mathbb{0}))\eta \subseteq \bigcup_{a \in N} \mathbb{E}(\theta_\alpha^{\rightsquigarrow}(a))\eta$

where  $\alpha_1, \alpha_2$  are fresh variables,  $\theta_\alpha^{\rightarrow}(a) = a\{(\alpha_1 \rightarrow \alpha_2) \vee \alpha / \alpha\}$ , and  $\theta_\alpha^{\rightsquigarrow}(a) = a\{((\alpha_1 \rightarrow \alpha_2) \setminus (\mathbb{1} \rightarrow \mathbb{0})) \vee \alpha / \alpha\}$ .

As an aside notice that both lemmas above would not hold if the variable  $\alpha$  in their statements occurred negated at toplevel, whence the necessity of *Step 4*.

Finally, *Step 6* is justified by the two following lemmas in whose proofs the hypothesis of convexity plays a crucial role:

**Lemma 3.11.** *Let  $(\mathcal{D}, \llbracket \cdot \rrbracket)$  be a convex set-theoretic interpretation and  $P, N$  two finite subsets of  $\mathcal{A}_{\text{prod}}$ . Then:*

$$\forall \eta . \bigcap_{a \in P} \mathbb{E}(a)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \Leftrightarrow \forall N' \subseteq N . \begin{cases} \forall \eta . \llbracket \bigwedge_{t_1 \times t_2 \in P} t_1 \wedge \bigwedge_{t_1 \times t_2 \in N'} \neg t_1 \rrbracket \eta = \emptyset \\ \vee \\ \forall \eta . \llbracket \bigwedge_{t_1 \times t_2 \in P} t_2 \wedge \bigwedge_{t_1 \times t_2 \in N \setminus N'} \neg t_2 \rrbracket \eta = \emptyset \end{cases}$$

(with the convention  $\bigcap_{a \in \emptyset} \mathbb{E}(a)\eta = \mathbb{E}^{\text{prod}} \mathcal{D}$ ).

**Lemma 3.12.** *Let  $(\mathcal{D}, \llbracket \cdot \rrbracket)$  be a convex set-theoretic interpretation and  $P, N$  be two finite subsets of  $\mathcal{A}_{\text{fun}}$ . Then:*

$$\forall \eta . \bigcap_{a \in P} \mathbb{E}(a)\eta \subseteq \bigcup_{a \in N} \mathbb{E}(a)\eta \Leftrightarrow \begin{cases} \forall \eta . \llbracket t_0 \setminus (\bigvee_{t \rightarrow s \in P'} t) \rrbracket \eta = \emptyset \\ \vee \\ \exists (t_0 \rightarrow s_0) \in N . \forall P' \subseteq P . \begin{cases} P \neq P' \\ \wedge \\ \forall \eta . \llbracket (\bigwedge_{t \rightarrow s \in P \setminus P'} s) \setminus s_0 \rrbracket \eta = \emptyset \end{cases} \end{cases}$$

(with the convention  $\bigcap_{a \in \emptyset} \mathbb{E}(a)\eta = \mathbb{E}^{\text{fun}} \mathcal{D}$ )

### 3.4 Algorithm

The formalization of the subtyping algorithm is done via notion of *simulation* that we borrow from [10] and extend to account for type variables and type instantiation. We use  $\theta$  to range over (syntactic) substitutions, that is, partial functions from  $\mathcal{V}$  to  $\mathcal{T}$  (we reserve the meta-variable  $\sigma$  we introduced in Section 2 for ground substitutions). The application of a substitution to a type is straightforwardly homomorphically defined (notice that there is no binder). This definition is naturally extended to normal forms by applying the substitution to each type in the sets that form the normal form. It is then used to define the set of instances of a type.

**Definition 3.13 (Instances).** *Given a type  $t \in \mathcal{T}$ , we define  $[t]_{\approx}$ , the set of instances of  $t$  as:*

$$[t]_{\approx} \stackrel{\text{def}}{=} \{s \mid \exists \theta : \mathcal{V} \rightarrow \mathcal{T} . t\theta = s\}$$

**Definition 3.14 (Simulation).** *Let  $\mathcal{S}$  be an arbitrary set of normal forms. We define another set of normal forms  $\mathbb{E}\mathcal{S}$  as*

$$\{\tau \mid \forall s \in [\tau]_{\approx} . \forall u \in U . \forall (P, N) \in \mathcal{N}(s) . (P \subseteq \mathcal{A}_u \Rightarrow C_u^{P, N \cap \mathcal{A}_u})\}$$

where:

$$C_{\text{basic}}^{P, N} \stackrel{\text{def}}{=} \bigcap_{b \in P} \mathbb{B}(b) \subseteq \bigcup_{b \in N \cap \mathcal{A}_{\text{basic}}} \mathbb{B}(b)$$

$$C_{\text{prod}}^{P, N} \stackrel{\text{def}}{=} \forall N' \subseteq N . \begin{cases} \mathcal{N}(\bigwedge_{t_1 \times t_2 \in P} t_1 \wedge \bigwedge_{t_1 \times t_2 \in N'} \neg t_1) \in \mathcal{S} \\ \vee \\ \mathcal{N}(\bigwedge_{t_1 \times t_2 \in P} t_2 \wedge \bigwedge_{t_1 \times t_2 \in N \setminus N'} \neg t_2) \in \mathcal{S} \end{cases}$$

$$C_{\text{fun}}^{P, N} \stackrel{\text{def}}{=} \exists t_0 \rightarrow s_0 \in N . \forall P' \subseteq P .$$

$$\begin{cases} \mathcal{N}(t_0 \wedge \bigwedge_{t \rightarrow s \in P'} \neg t) \in \mathcal{S} \\ \vee \\ \begin{cases} P \neq P' \\ \wedge \\ \mathcal{N}(\neg s_0 \wedge \bigwedge_{t \rightarrow s \in P \setminus P'} s) \in \mathcal{S} \end{cases} \end{cases}$$

We say that  $\mathcal{S}$  is a simulation if:

$$\mathcal{S} \subseteq \mathbb{E}\mathcal{S}$$

The notion of simulation is at the basis of our subtyping algorithm. In what follows we show that simulations soundly and completely characterize the set of empty types of a well-founded convex model. More precisely, we show that every type in a simulation is empty (soundness) and that the set of all empty types is a simulation (completeness), actually, the largest simulation.

**Lemma 3.15.** *Let  $\llbracket \cdot \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a set-theoretic interpretation and  $t$  a type. If  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \llbracket t \rrbracket \eta = \emptyset$ , then  $\forall \theta : \mathcal{V} \rightarrow \mathcal{T} . \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \llbracket t\theta \rrbracket \eta = \emptyset$*

Lemma 3.15 shows that if a type is empty, then all its syntactic instances are empty. In particular if a type is empty we can rename all its type variables without changing any property. So when working with empty types or, equivalently, with subtyping relations, types can be considered equivalent modulo  $\alpha$ -renaming (ie, the renaming of type variables).

The first result we prove is that every simulation contains only empty types.

**Theorem 3.16 (Soundness).** *Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a convex structural interpretation and  $\mathcal{S}$  a simulation. Then for all  $\tau \in \mathcal{S}$ , we have  $\forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \llbracket \tau \rrbracket \eta = \emptyset$*

The intuition of the simulation is that if we consider the statements of Lemmas 3.11 and 3.12 as if they were rewriting rules (from right to left), then  $\mathbb{E}\mathcal{S}$  contains all the types that we can deduce to be empty in one step reduction when we suppose that the types in  $\mathcal{S}$  are empty. A simulation is thus a set that is already saturated with respect to such a rewriting. In particular, if we consider the statements of Lemmas 3.11 and 3.12 as inference rules for determining when a type is equal to  $\emptyset$ , then  $\mathbb{E}\mathcal{S}$  is the set of immediate consequences of  $\mathcal{S}$ , and a simulation is a *self-justifying* set, that is a co-inductive proof of the fact that all its elements are equal to  $\emptyset$ .

Completeness derives straightforwardly from the construction of simulation and the lemmas we proved about it.

**Theorem 3.17.** *Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a convex set-theoretic interpretation. We define a set of normal forms  $\mathcal{S}$  by:*

$$\mathcal{S} \stackrel{\text{def}}{=} \{ \tau \mid \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \llbracket \tau \rrbracket \eta = \emptyset \}$$

Then:

$$\mathbb{E}\mathcal{S} = \{ \tau \mid \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \mathbb{E}(\tau)\eta = \emptyset \}$$

**Corollary 3.18.** *Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a convex structural interpretation. Define as above  $\mathcal{S} = \{ \tau \mid \forall \eta \in \mathcal{P}(\mathcal{D})^\vee . \mathbb{E}(\tau)\eta = \emptyset \}$ . Then  $\llbracket \_ \rrbracket$  is a model if and only if  $\mathcal{S} = \mathbb{E}\mathcal{S}$ .*

This corollary has two implications: first, that the condition for a set-theoretic interpretation to be a model depends only on the subtyping relation it induces; second, that the simulation that contains all the empty types is the largest simulation. In particular this second implication entails the following corollary.

**Corollary 3.19 (Completeness).** *Let  $\llbracket \_ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{D})^\vee \rightarrow \mathcal{P}(\mathcal{D})$  be a well-formed convex model and  $s$  and  $t$  two types. Then  $s \leq t$  if and only if there exists a simulation  $\mathcal{S}$  such that  $\mathcal{N}(s \wedge \neg t) \in \mathcal{S}$ .*

The corollary states that to check whether  $s \leq t$  we have to check whether there exists a simulation that contains the normal form, denoted by  $\tau_0$ , of  $s \wedge \neg t$ . Thus a simple subtyping algorithm is to use Definition 3.14 as a set of saturation rules: we start from  $\tau_0$  and try to saturate it. At each step of saturation we add just the normal forms in which we eliminated top-level variables according to Lemmas 3.8, 3.9, and 3.10. Because of the presence of “or” in the definition, the algorithm follows different branches until it reaches a simulation (in which case it stops with success) or it adds a non-empty type (in which case the whole branch is abandoned).

All results we stated so far have never used the regularity of types. The theory holds also for non regular types and the algorithm described above is a sound and complete procedure to check their inclusion. The only result that needs regularity is decidability.

### 3.5 Decidability

As anticipated in the related work section, the subtyping relation on polymorphic regular types is decidable in EXPTIME. We prove

decidability but the result on complexity is due to Gesbert *et al.* [14] who, as we already hinted, gave a linear encoding of the relation presented here in their variant of the  $\mu$ -calculus, for which they have an EXPTIME solver [13], thus obtaining a subtyping decision algorithm that is EXPTIME<sup>8</sup> (in doing so they also spotted a subtle error in the original definition of our subtyping algorithm).

To prove decidability we just prove that our algorithm terminates. We do so by first showing that it terminates on finite trees (which is the crux of the problem since the only potential source of loop are the substitutions of free variables performed in **Step 5**) and then show that when we switch to regular trees the algorithm always ends up on memoized terms (see the extended version).

### 3.6 Convex Models

The last step of our development is to prove that there exists at least one set-theoretical model that is convex. From a practical point of view this step is not necessary since one can always take the work we did so far as a syntactic definition of the subtyping relation. However, the mathematical support makes our theory general and applicable to several different domains (eg, Gesbert *et al.*'s starting point and early attempts relied on this result), so finding a model is not done just for the sake of the theory. As it turns out, there actually exist a lot of convex models since every model for ground types with infinite denotations is convex. So to define a convex model it just suffices to take any model defined in [10] and straightforwardly modify the interpretation of basic and singleton types (more generally, of all indivisible types<sup>9</sup>) so they have infinite denotations.

**Definition 3.20 (Infinite support).** *A model  $(\mathcal{D}, \llbracket \_ \rrbracket)$  is with infinite support if for every ground type  $t$  and assignment  $\eta$ , if  $\llbracket t \rrbracket \eta \neq \emptyset$ , then  $\llbracket t \rrbracket \eta$  is an infinite set.*

What we want to prove then is the following theorem.

**Theorem 3.21.** *Every well-founded model with infinite support is convex.*

The proof of this theorem is quite technical—it is the proof that required us most effort—and proceeds in three logical steps. First, we prove that the theorem holds when all types at issue do not contain any product or arrow type. In other words, we prove that equation (8) holds for  $\llbracket \_ \rrbracket$  with infinite support and where all  $t_i$ 's are Boolean combinations of type variables and basic types. This is the key step in which we use the hypothesis of infinite support, since in the proof—done by contradiction—we need to pick an unbounded number of elements from our basic types in order to build a counterexample that proves the result. Second, we extend the previous proof to any type  $t_i$  that contains finitely many applications of the product constructor. In other terms, we prove the result for any (possibly infinite) type, provided that recursion traverses just arrow types, but not product types. As in the first case, the proof builds some particular elements of the domain. In the presence of type constructors the elements are built inductively. This is possible since products are not recursive, while for arrows it is always possible to pick a fresh appropriate element that resides in that arrow since every arrow type contains the (indivisible) closed type  $\mathbb{1} \rightarrow \emptyset$ . Third, and last, we use the previous result and the well-foundedness of the model to show that the result holds for all types, that is, also for types in which recursion traverses a product type. More precisely, we prove that if we assume that the result

<sup>8</sup>This is also a lower bound for the complexity since the subtyping problem for ground regular types (without arrows) is known to be EXPTIME-complete.

<sup>9</sup>Our system has a very peculiar indivisible type:  $\mathbb{1} \rightarrow \emptyset$ , the type of the functions that diverge on all arguments. This can be handled by adding a fixed infinite set of fresh elements of the domain to the interpretation of every arrow type (cf. the proof of Corollary 3.26 in the extended version).

does not hold for some infinite product type then it is possible to build a finite product type (actually, a finite expansion of the infinite type) that disproves equation (8), contradicting what is stated in our second step. Well-foundedness of the model allows us to build this finite type by induction on the elements denoted by the infinite one.

More precisely, we proceed as follows.

**Definition 3.22.** We use  $\mathcal{T}^{\text{fp}}$  to denote the set of types with finite products, that is, the set of all types in which every infinite branch contains a finite number of occurrences of the  $\times$  constructor.

The first two steps of our proof are proved simultaneously in the following lemma.

**Lemma 3.23.** Let  $(\mathcal{D}, \llbracket - \rrbracket)$  be a model with infinite support, and  $t_i \in \mathcal{T}^{\text{fp}}$  for  $i \in [1..n]$ . Then

$$\begin{aligned} \forall \eta . \llbracket t_1 \rrbracket \eta = \emptyset \vee \dots \vee \llbracket t_n \rrbracket \eta = \emptyset &\Leftrightarrow \\ \forall \eta . \llbracket t_1 \rrbracket \eta = \emptyset \vee \dots \vee \forall \eta . \llbracket t_n \rrbracket \eta = \emptyset & \end{aligned}$$

Finally, it just remains to prove Theorem 3.21, that is to say, that Lemma 3.23 above holds also for  $t_i$ 's with recursive products. This result requires the following preliminary lemma.

**Lemma 3.24.** Let  $\llbracket - \rrbracket$  be a well-founded model with infinite support and  $t$  a type (which may thus contain infinite product types). If there exists a value  $d$  and an assignment  $\bar{\eta}$  such that  $d \in \llbracket t \rrbracket \bar{\eta}$ , then there exists a type  $t^{\text{fp}} \in \mathcal{T}^{\text{fp}}$  such that  $d \in \llbracket t^{\text{fp}} \rrbracket \bar{\eta}$  and for all assignment  $\eta$  if  $\llbracket t \rrbracket \eta = \emptyset$ , then  $\llbracket t^{\text{fp}} \rrbracket \eta = \emptyset$ .

While the statement of the previous lemma may, at first sight, seem obscure, its meaning is rather obvious. It states that in a well-founded model (ie, a model in which all the values are finite) whenever a recursive (product) type contains some value, then we can find a *finite* expansion of this type that contains the same value; furthermore, if the recursive type is empty in a given assignment, then also its finite expansion is empty in that assignment. This immediately yields our final result.

**Lemma 3.25.** Let  $(\mathcal{D}, \llbracket - \rrbracket)$  be a well-founded model with infinite support, and  $t_i$  for  $i \in [1..n]$ . Then

$$\begin{aligned} \forall \eta . \llbracket t_1 \rrbracket \eta = \emptyset \vee \dots \vee \llbracket t_n \rrbracket \eta = \emptyset &\Leftrightarrow \\ \forall \eta . \llbracket t_1 \rrbracket \eta = \emptyset \vee \dots \vee \forall \eta . \llbracket t_n \rrbracket \eta = \emptyset & \end{aligned}$$

**Corollary 3.26 (Convex model).** *There exists a convex model.*

## 4. Conclusion

This work constitutes the first solution to the problem of defining a semantic subtyping relation for a polymorphic extension of regular tree types. This problem, despite its important practical interest and potential fallouts, has been somehow neglected by most recent research since it was considered untreatable or unfeasible. Our solution not only has immediate application to the definition of programming languages for XML, but also opens several new research directions that we briefly discuss below.

The first direction concerns the definition of extensions of the type system itself. Among the possible extensions the most interesting (and difficult) one seems to be the extension of types with explicit second order quantifications. Currently, we consider prenex polymorphism, thus quantification on types is performed only at meta-level. But since this work proved the feasibility of a semantic subtyping approach for polymorphic types, we are eager to check whether it can be further extended to impredicative second order types, by adding explicit type quantification. This would be interesting not only from a programming language perspective, but also from a logic viewpoint since it would remove some of the limitations to the introspection capabilities we pointed out in Section 2.7. This may move our type system closer to being an expres-

sive logic for subtyping. On the model side, it would be interesting to check whether the infinite support property (Definition 3.20) is not only a sufficient but also a necessary condition for convexity. This seems likely to the extent that the result holds for the type system restricted to basic type constructors (ie, without products and arrows). However, this is just a weak conjecture since the proof of sufficiency heavily relies (in the case for product types) on the well-foundedness property. Therefore, there may even exist non-well-founded models (non-well-founded models exist in the ground case: see [8, 10]) that are with infinite support but not convex. Nevertheless, an equivalent characterization of convexity—whether it were infinite support or some other characterization—would provide us a different angle of attack to study the connections between convexity and parametricity (see later on).

The second direction is to explore the definition of new languages to take advantage of the new capabilities of our system. A first natural test will be to see how to add overloaded (typed by intersection types) and higher-order (typed by arrow types) functions to the language defined in [15]. This already looks as quite a challenging problem since it needs local type inference for both subtyping and instantiation,<sup>10</sup> and we are actively working on it. But the overall design space for a programming language that can exploit the advanced features of our types is rather large since a lot of possible variations can be considered (eg, the use of type variables in pattern matching) and even more features can be encoded (eg, as explained in Footnote 4, bounded quantification can not only be encoded via intersection types but, thanks to them, also made more general since intersections allow the programmer to specify bounds on a per-occurrence basis). While exploring this design space it will be interesting to check whether our polymorphic union types can encode advanced type features such as polymorphic variants [12] and GADTs [24].

In our opinion, the definition of convexity is the most important and promising contribution of our work especially in view of its potential implications on the study of parametricity. As a matter of fact, there are strong connections between parametricity and convexity. We have already seen that convexity removes the stuttering behavior that clashes with parametricity, as equation (5) clearly illustrates. More generally, both convexity and parametricity describe or enforce uniformity of behavior. Parametricity imposes to functions a uniform behavior on parameters typed by type variables, since the latter cannot be deconstructed. This allows Wadler to deduce “theorems for free”: the uniformity imposed by parametricity (actually, imposed by the property of being definable in second order  $\lambda$ -calculus) dramatically restricts the choice of possible behaviors of parametric functions to a point that it is easy to deduce theorems about a function just by considering its type [22]. In a similar way convexity imposes a uniform behavior to the zeros of the semantic interpretation, which is equivalent to imposing uniformity to the subtyping relation. An example of this uniformity is given by product types: in our framework a product  $(t_1 \times \dots \times t_n)$  is empty (ie, it is empty for every possible assignment of its variables) if and only if there exists a particular  $t_i$  that is empty (for

<sup>10</sup>To have some flavor of the problem, consider an overloaded function even for the domain  $\mathbb{1}$  that when applied to an integer returns whether it is even or not, while it returns arguments of any other type unchanged. Its type is  $(\text{Int} \rightarrow \text{Bool}) \wedge (\alpha \setminus \text{Int} \rightarrow \alpha \setminus \text{Int})$ . If we apply a curried  $\text{map} : (\beta \rightarrow \gamma) \rightarrow \beta^* \rightarrow \gamma^*$  to  $\text{odd}$ , then we want to deduce  $\text{map}(\text{even}) : (\text{Int}^* \rightarrow \text{Bool}^*) \wedge ((\alpha \setminus \text{Int})^* \rightarrow (\alpha \setminus \text{Int})^*) \wedge ((\alpha \vee \text{Int})^* \rightarrow ((\alpha \setminus \text{Int}) \vee \text{Bool})^*)$ . That is, when we apply  $\text{map}(\text{even})$  to a list of integers it returns a list of booleans; when we apply it to a list that does not contain integers, then it returns a list of the same type; when the list contains some integers (eg, a list of reals), then it replaces integer elements by boolean ones. This is not an instance of the output type of  $\text{map}$ .

all possible assignments). We recover a property typical of closed types.

We conjecture the connection to be much deeper than described above. This can be clearly perceived by rereading the original Reynolds paper on parametricity [18] in the light of our results. Reynolds tries to characterize parametricity—or *abstraction* in Reynolds terminology—in a set-theoretic setting since, in Reynolds words, “if types denote specific subsets of a universe then their unions and intersections are well defined”, which in Reynolds opinion is the very idea of abstraction. This can be rephrased as the fact that the operations for some types are well defined independently from the representation used for each type (Reynolds speaks of abstraction and representation since he sees the abstraction property as a result about change of representation). The underlying idea of parametricity according to Reynolds is that “meanings of an expression in ‘related’ environments will be ‘related’” [18]. But as he points out few lines later “while the relation for each type variable is arbitrary, the relation for compound type expressions [*ie*, type constructors] must be induced in a specified way. We must specify how an assignment of relations to type variables is extended to type expressions” [18]. Reynolds formalizes this extension by defining a “relation semantics” for type constructors and, as Wadler brilliantly explains [22], this corresponds to regard types as relations. In particular pairs are related if their corresponding components are related and functions are related if they take related arguments into related results: there is a precise correspondence with the extensional interpretation of type constructors we gave in Definition 3.5 and, more generally, between the framework used to state parametricity and the one in our work.

Parametricity holds for terms written in the Girard/Reynolds second order typed lambda calculus (also known as pure polymorphic lambda calculus or System F). The property of being definable in the second-order typed lambda-calculus is *the* condition that harnesses expressions and forces them to behave uniformly. Convexity, instead, does not require any definability property. It *semantically* harnesses the denotations of expressions and forces them to behave uniformly. Therefore all seems to suggest that convexity may be a semantic characterization of what in Reynolds approach is the definability in the second-order typed lambda-calculus, which is a syntactic property. Or, to put it otherwise, convexity states parametricity for (or transposes it to) models rather than languages.

Although we have this strong intuition about the connection between convexity and parametricity, we do not know how to express this connection in a formal way, yet. We believe that the answer may come from the study of the calculus associated to our subtyping relation. We do not speak here of some language that can take advantage of our subtyping relation and whose design space we discussed earlier in this conclusion. What we are speaking of is every calculus whose model of values (*ie*, the model obtained by associating each type to the set of values that have that type) induces the same subtyping relation as the one devised here. Indeed, as parametricity leaves little freedom to the definition of transformations, so the semantic subtyping framework leaves little freedom to the definition of a language whose model of values induces the same subtyping relation as the relation used to type its values. If we could prove that every such language must automatically satisfy Reynolds abstraction theorem (and, even more, prove also the converse), then we would have a formal and strong connection between convexity and parametricity, the former being a purely semantic (in the sense that it does not rely on any language or calculus) characterization of the latter. But this is a long term and ambitious goal that goes well beyond the scope of the present work.

**Acknowledgments.** We were given a decisive help by Nino Salibra, who showed us how to prove that all powersets of infinite sets are convex models of the propositional logic. Pietro Abate implemented an early version of our prototype subtype checker. Bow-Yaw Wang drew our attention to the connection with convex theories. Claude Benzaken pointed out several useful results in combinatorics to us. Nils Gesbert spotted a subtle mistake in our subtyping algorithm. Many persons discussed the contents of this work with us and gave precious suggestions. Among them we are particularly grateful to Véronique Benzaken, Mariangiola Dezani-Ciancaglini, Kim Nguyen, Daniele Varacca, and Philip Wadler.

## References

- [1] V. Balat, J. Vouillon, and B. Yakobowski. Experience report: Ocsigen, a web programming framework. In *ICFP '09*. ACM Press, 2009.
- [2] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-friendly general purpose language. In *ICFP '03*. ACM Press, 2003.
- [3] G. Castagna, R. De Nicola, and D. Varacca. Semantic subtyping for the  $\pi$ -calculus. *Theor. Comput. Sci.*, 398(1-3):217–242, 2008.
- [4] G. Castagna and K. Nguyen. Typed iterators for XML. In *ICFP '08*, pages 15–26. ACM Press, 2008.
- [5] J. Clark and M. Murata. Relax-NG, 2001. [www.relaxng.org](http://www.relaxng.org).
- [6] D. Draper et al. XQuery 1.0 and XPath 2.0 Formal Semantics, 2007. <http://www.w3.org/TR/query-semantics/>.
- [7] T. Berners-Lee et al. *Uniform Resource Identifier*, January 2005. RFC 3986, STD 66.
- [8] A. Frisch. *Théorie, conception et réalisation d'un langage de programmation fonctionnel adapté à XML*. PhD thesis, Université Paris 7, December 2004.
- [9] A. Frisch. OCaml + XDuce. In *ICFP '06*. ACM Press, 2006.
- [10] A. Frisch, G. Castagna, and V. Benzaken. Semantic subtyping: dealing set-theoretically with function, union, intersection, and negation types. *The Journal of ACM*, 55(4):1–64, 2008.
- [11] V. Gapeyev, M.Y. Levin, B.C. Pierce, and A. Schmitt. The Xtatic compiler and runtime system, 2005. <http://www.cis.upenn.edu/~bcpierce/xtatic>.
- [12] J. Garrigue. Programming with polymorphic variants. In *Proc. of ML Workshop*, 1998.
- [13] P. Genevès, N. Layaïda, and A. Schmitt. Efficient static analysis of XML paths and types. In *PLDI '07*. ACM Press, 2007.
- [14] N. Gesbert, P. Genevès, and N. Layaïda. Parametric Polymorphism and Semantic Subtyping: the Logical Connection. In *ICFP '11*, 2011.
- [15] H. Hosoya, A. Frisch, and G. Castagna. Parametric polymorphism for XML. *ACM TOPLAS*, 32(1):1–56, 2009.
- [16] H. Hosoya and B. Pierce. XDuce: A statically typed XML processing language. *ACM TOIT*, 3(2):117–148, 2003.
- [17] K. Zhuo Ming Lu and M. Sulzmann. An implementation of subtyping among regular expression types. In *Proc. of APLAS'04*, volume 3302 of *LNCS*, pages 57–73. Springer, 2004.
- [18] J.C. Reynolds. Types, abstraction and parametric polymorphism. In *Information Processing*, pages 513–523. Elsevier, 1983.
- [19] J.C. Reynolds. Polymorphism is not set-theoretic. In *Semantics of Data Types*, volume 173 of *LNCS*, pages 145–156. Springer, 1984.
- [20] J. Vouillon. Polymorphic regular tree types and patterns. In *POPL '06*, pages 103–114, 2006.
- [21] W3C. *SOAP Version 1.2*. <http://www.w3.org/TR/soap>.
- [22] P. Wadler. Theorems for free! In *Functional Programming Languages and Computer Architecture*, pages 347–359. ACM Press, 1989.
- [23] C. Wallace and C. Runciman. Haskell and XML: Generic combinators or type based translation? In *ICFP '99*, pages 148–159, 1999.
- [24] H. Xi, C. Chen, and G. Chen. Guarded recursive datatype constructors. In *POPL '03*, pages 224–235. ACM Press, 2003.