# A Constraint-Based Presentation and Generalization of Rows

François Pottier

INRIA

## Abstract

*We study the combination of possibly conditional non-structural subtyping constraints with rows. We give a new presentation of rows, where row terms disappear; instead, we annotate constraints with* filters. *We argue that, in the presence of subtyping, this approach is simpler and more general. In the case where filters are finite or cofinite sets of row labels, we give a constraint solving algorithm whose complexity is $O(n^3 m \log m)$, where $n$ is the size of the constraint and $m$ is the number of row labels that appear in it. We point out that this allows efficient type inference for record concatenation. Furthermore, by varying the nature of filters, we obtain several natural generalizations of rows.*

## 1 Introduction

*Records* are an important feature of programming languages, not only because of their ubiquitous use in defining data structures, but also because they form a basic layer on top of which other features, such as objects and modules, may be built. Records are finite associations of values to labels; typical primitive operations on records are access (extracting the contents of a field), update (modifying the contents of an existing field), extension (adding a new field), restriction (removing an existing field), and concatenation (merging two existing records into a new one) [1].

*Rows* [16, 14] offer syntax to describe infinite families of types, indexed by labels. Rows offer random access to any single component. Furthermore, they offer uniform access to *every* component at once; that is, by imposing a constraint on rows, one effectively imposes a family of constraints, point-wise, on their components. Equality constraints on rows allow assigning accurate types to record access, update, extension, and restriction.

*Subtyping* is a common way of increasing a type system's expressiveness. In the presence of record types, subtyping is usually *non-structural*—that is, a type and its subtypes may not have the same shape—because it allows forgetting whether certain fields are present in a record. In the presence of rows, this phenomenon occurs at the level

of field types. Indeed, typical field types are Abs (the field does not exist in this record), Pre $\tau$ (the field exists and contains a value of type $\tau$) and $\top_{field}$ (the field may or may not exist). The definition of subtyping typically involves axioms such as Pre $\tau \leq \top_{field}$ and Abs $\leq \top_{field}$, the first of which is non-structural.

From a type-theoretic perspective, concatenation is the most challenging operation on records, because the type of its result cannot be related to the types of its arguments using subtyping constraints only. Consider, for instance, *symmetric* concatenation, which requires its arguments to have disjoint sets of fields. For every field $\ell$, the result type at $\ell$ is Pre $\tau$ if one argument has type Pre $\tau$ at $\ell$ and the other has type Abs at $\ell$; it is Abs if both arguments have type Abs at $\ell$; the operation is ill-typed otherwise. The difficulty is not in the quantification over $\ell$, which is dealt with by rows, but in the fact that this definition is *by cases*. In [11], I suggested addressing this issue using *conditional constraints*, a concept whose origin can be traced back to Reynolds [15].

Non-structural subtyping constraints can be solved in cubic time [6]. In fact, McAllester's theorem [4] allows establishing that possibly conditional subtyping constraints can be solved in cubic time as well. However, the combination of such constraints with rows, which is implemented in [10], has never been studied from a complexity-theoretic point of view. Taking advantage of this fact, Palsberg and Zhao [7] recently claimed to have developed the first polynomial time type inference algorithm for symmetric record concatenation. Their algorithm runs in time $O(n^5)$, where $n$ is the size of the program. It is not compositional, that is, it must analyze a whole program at once, which makes it rather impractical.

The contribution of this paper is two-fold. First, we give a new presentation of rows, where row *terms* disappear; instead, we annotate *constraints* with *filters*, which, in the simplest case, are finite or cofinite sets of row labels. We argue that this presentation is simpler and more general. Second, we give a solver for this new constraint language, and establish that it runs in time $O(n^3 m \log m)$, where $n$ is the size of the constraint and $m$ is the number of row labels that appear in it. As an application, we show that type inference for symmetric record concatenation can be reduced to solv-

ing a constraint whose size is linear in the program's size. Thus, this approach outperforms Palsberg and Zhao's. It is, furthermore, compositional.

Why a new account of rows? In existing presentations, rows are *terms*. In particular, if $\ell$ is a row label, $\tau$ is a type and $\rho$ is a row, then the term $(\ell : \tau; \rho)$ is a row as well. In type systems *without* subtyping, this allows reducing type inference to first-order unification in an equational theory, whose axioms allow e.g. commuting the order of row labels in a row term. This is crucial, because unification is extremely efficient. One disadvantage, however, is that row unification sometimes allocates fresh type variables, making termination and complexity arguments more delicate. Type systems *with* subtyping, on the other hand, replace unification with a more costly constraint solving procedure. The key insight, then, is that there is no longer a point in using row *terms*. One may instead annotate *constraints* with *filters*, which carry information about row labels. This has two consequences. First, it becomes possible to solve constraints without allocating fresh type variables. This simplifies the complexity analysis, and was our original motivation. Second, even though we initially define filters as finite or cofinite sets of row labels, most of our development (in fact, all of it but the complexity analysis) is independent of the nature of filters, as long as they are preserved by Boolean operations and enjoy a decidable emptiness test. Thus, by varying the nature of filters, we obtain several natural generalizations of rows.

The paper is laid out as follows. Section 2 presents the syntax, as well as the meaning, of types and constraints. Section 3 describes a satisfiability test for constraints, based on a closure computation, and assesses its theoretical complexity. Section 4 explains how type inference for record operations, including symmetric concatenation, reduces to constraint solving. Section 5 discusses a few extensions of the framework.

## 2  Formal Presentation of the System

In order to achieve a good measure of generality, our presentation is parameterized with respect to a *ground signature*, which specifies a type algebra and a subtype ordering; it is defined in Section 2.1. Given a fixed ground signature, Section 2.2 defines the algebra of ground terms and its ordering. The syntax of types and constraints is given in Section 2.3 and their logical interpretation in Section 2.4.

### 2.1  Assumptions

**Definition 1** *Let a* sort *$\varsigma$ be one of* $\{Type, Row\}$.

A ground signature consists of three components: a family of *symbol* lattices, indexed by *kinds*, a set of *parameter labels*, each of which is co- or contra-variant and has fixed sort and kind, and a description of each symbol's *arity* as a set of parameter labels.

**Definition 2** *Let $\mathcal{K}$ be a finite set of* kinds*. For every kind $\kappa \in \mathcal{K}$, let $\mathcal{S}_\kappa$ be a finite lattice of* symbols*, with operations $\perp_\kappa$, $\top_\kappa$, $\leq_\kappa$, $\sqcup_\kappa$ and $\sqcap_\kappa$. We write $\mathcal{S}$ for the set of all symbols. Let $\mathcal{P}^+$ and $\mathcal{P}^-$ be disjoint finite sets of* parameter labels*. We write $\mathcal{P}$ for $\mathcal{P}^+ \uplus \mathcal{P}^-$. With every parameter label $p \in \mathcal{P}$, associate a sort and a kind, written* sort$(p)$ *and* kind$(p)$*, respectively. With every symbol $s \in \mathcal{S}$, associate a subset of $\mathcal{P}$, called the* arity *of $s$ and written $a(s)$. We require that, for all $\kappa \in \mathcal{K}$ and $s_0, s_1, s_2 \in \mathcal{S}_\kappa$, $s_0 \leq_\kappa s_1 \leq_\kappa s_2$ imply $a(s_0) \cap a(s_2) \subseteq a(s_1)$.*

The last requirement is used in the proof of Theorem 1 to ensure that the orderings $\leq_\kappa$ do give rise to an ordering on ground terms. The information described in Definition 2 forms a *ground signature*. We now assume that some ground signature is fixed.

**Example 1** Let $\mathcal{K}$ be $\{type, field\}$. Let $\mathcal{S}_{type}$ be the flat lattice whose elements other than $\perp_{type}$ and $\top_{type}$ are $\rightarrow$ and $\{\cdot\}$. Let $\mathcal{S}_{field}$ be the flat lattice whose elements other than $\perp_{field}$ and $\top_{field}$ are Abs and Pre. Let $\mathcal{P}^- = \{dom\}$ and $\mathcal{P}^+ = \{rng, fields, content\}$. Let kind$(dom) =$ kind$(rng) =$ kind$(content) = type$ and kind$(fields) = field$. Let sort$(dom) =$ sort$(rng) =$ sort$(content) = Type$ and sort$(fields) = Row$. Let $a(\rightarrow) = \{dom, rng\}$ and $a(\{\cdot\}) = \{fields\}$ and $a(\mathsf{Pre}) = \{content\}$. Let the arity of all other symbols be $\varnothing$.

This ground signature defines two kinds of types, plain types and field types. The record type constructor $\{\cdot\}$ forms a plain type out of a row of field types. The constructor Pre forms a field type out of a plain type. Abs and Pre are made incomparable, because that is required to assign a sound type to record concatenation [1, 11], but they do have a common supertype $\top_{field}$, so *width subtyping* is present.

If we take a look ahead at the syntax of terms and at the sorting and kinding restrictions, defined in Section 2.3, we find that this ground signature gives rise to the following grammar of terms, where $\tau$ ranges over terms of sort *Type* and kind *type*, $\rho$ ranges over terms of sort *Row* and kind *field*, and $\varphi$ ranges over terms of sort *Type* and kind *field*:

$$\tau ::= \alpha \mid \perp_{type} \mid \top_{type} \mid \tau \rightarrow \tau \mid \{\rho\}$$
$$\rho ::= \alpha \mid \partial\varphi$$
$$\varphi ::= \alpha \mid \perp_{field} \mid \top_{field} \mid \mathsf{Pre}\,\tau \mid \mathsf{Abs}$$

This concrete style of definition is standard [14, 11]. Parameterizing our development with respect to a ground signature allows us to accommodate a whole family of definitions in this style.

## 2.2 Logical Model

We now define the universe of ground terms, which forms the model within which types and constraints are interpreted. As in [12, 7], we allow types to be recursive, that is, we define ground terms as infinite trees. We do not restrict our attention to regular trees, because that would not affect constraint satisfiability or entailment. Since the model contains infinite trees, its definition is naturally coinductive. We cannot devote space to establishing the existence of the fixed points mentioned below; a detailed treatment of a similar construction can be found in [3].

Let $\mathcal{L}$ be a denumerable set of *row labels*. Let $\ell$ range over row labels.

For the sake of brevity, let us write $s(\bar{t}_p)$ for $s(t_p)_{p\in a(s)}$ and $(\bar{t}_\ell)$ for $(t_\ell)_{\ell\in\mathcal{L}}$. The former is a ground type whose head symbol is $s$ and whose sub-terms are the $t_p$, where $p$ ranges over the arity of $s$. The latter is a ground row, that is, a family of ground types indexed by row labels. This is made precise by the following definition:

**Definition 3** *The family of models $\mathbb{T}_\kappa^\varsigma$ is the greatest solution to the following equations:*

$$\mathbb{T}_\kappa^{Type} = \{s(\bar{t}_p) \;;\; s \in \mathcal{S}_\kappa \wedge \forall p \in a(s) \quad t_p \in \mathbb{T}_{\mathrm{kind}(p)}^{\mathrm{sort}(p)}\}$$
$$\mathbb{T}_\kappa^{Row} = \{(\bar{t}_\ell) \;;\; \forall \ell \in \mathcal{L} \quad t_\ell \in \mathbb{T}_\kappa^{Type}\}$$

*An element $t$ of $\mathbb{T}_\kappa^\varsigma$ is a* ground term *of sort $\varsigma$ and kind $\kappa$.*

Two ground types are comparable if and only if they have comparable head symbols and their sub-terms at every common parameter label $p$ are comparable as well—in the reverse direction, if $p$ is contravariant. Ground rows are compared point-wise and covariantly.

**Definition 4** *Let $t \,{}^p{\leq}_\kappa^\varsigma\, t'$ stand for $t \leq_\kappa^\varsigma t'$ if $p$ is a member of $\mathcal{P}^+$, and for $t' \leq_\kappa^\varsigma t$ otherwise. The family of relations $\leq_\kappa^\varsigma$ is the greatest solution to the following equivalences:*

$$s(\bar{t}_p) \leq_\kappa^{Type} s'(\bar{t}'_p)$$
$$\iff s \leq_\kappa s' \wedge \forall p \in a(s) \cap a(s') \quad t_p \,{}^p{\leq}_{\mathrm{kind}(p)}^{\mathrm{sort}(p)}\, t'_p$$

$$(\bar{t}_\ell) \leq_\kappa^{Row} (\bar{t}'_\ell)$$
$$\iff \forall \ell \in \mathcal{L} \quad t_\ell \leq_\kappa^{Type} t'_\ell$$

This definition may be made more explicit by viewing $\leq_\kappa^\varsigma$ as the intersection of the denumerable family $({}_k{\leq}_\kappa^\varsigma)_{k\geq 0}$, where every ${}_0{\leq}_\kappa^\varsigma$ is the full binary relation and ${}_{k+1}{\leq}_\kappa^\varsigma$ is given by

$$s(\bar{t}_p) \,{}_{k+1}{\leq}_\kappa^{Type}\, s'(\bar{t}'_p)$$
$$\iff s \leq_\kappa s' \wedge \forall p \in a(s) \cap a(s') \quad t_p \,{}_k^p{\leq}_{\mathrm{kind}(p)}^{\mathrm{sort}(p)}\, t'_p$$

$$(\bar{t}_\ell) \,{}_k{\leq}_\kappa^{Row}\, (\bar{t}'_\ell)$$
$$\iff \forall \ell \in \mathcal{L} \quad t_\ell \,{}_k{\leq}_\kappa^{Type}\, t'_\ell$$

**Theorem 1** *Every $(\mathbb{T}_\kappa^\varsigma, \leq_\kappa^\varsigma)$ forms a lattice.*

*Proof.* It is straightforward to establish, by induction over $k$, that every ${}_k{\leq}_\kappa^\varsigma$ is a preorder. In particular, establishing that ${}_{k+1}{\leq}_\kappa^{Type}$ is transitive requires exploiting the last requirement in Definition 2. Because preorders are preserved by intersection, every $\leq_\kappa^\varsigma$ is a preorder.

Let $=_\kappa^\varsigma$ be the kernel of the preorder $\leq_\kappa^\varsigma$. Since every $\leq_\kappa$ is antisymmetric, $=_\kappa^\varsigma$ satisfies the following equivalences:

$$s(\bar{t}_p) =_\kappa^{Type} s'(\bar{t}'_p)$$
$$\iff s = s' \wedge \forall p \in a(s) \quad t_p =_{\mathrm{kind}(p)}^{\mathrm{sort}(p)} t'_p$$

$$(\bar{t}_\ell) =_\kappa^{Row} (\bar{t}'_\ell)$$
$$\iff \forall \ell \in \mathcal{L} \quad t_\ell =_\kappa^{Type} t'_\ell$$

Because equality is precisely the greatest solution of these equivalences, $=_\kappa^\varsigma$ is in fact equality. So, $\leq_\kappa^\varsigma$ is an ordering.

Let $\square$ range over $\{\sqcup, \sqcap\}$. Let ${}^p\sqcup$ stand for $\sqcup$ if $p$ is a member of $\mathcal{P}^+$ and for $\sqcap$ otherwise. Define ${}^p\sqcap$ analogously. Then, the operations $\sqcup_\kappa^\varsigma$ and $\sqcap_\kappa^\varsigma$, whose argument is a finite subset of $\mathbb{T}_\kappa^\varsigma$, are defined coinductively by

$$\square_\kappa^{Type} \{s_i(\bar{t}_p^i)\}_{i\in I}$$
$$= (\square_\kappa \{s_i\}_{i\in I})\overline{\left({}^p\square_{\mathrm{kind}(p)}^{\mathrm{sort}(p)} \{t_p^i\}_{i\in I \wedge p\in a(s_i)}\right)}$$

$$\square_\kappa^{Row} \{(\bar{t}_\ell^i)\}_{i\in I}$$
$$= \overline{(\square_\kappa^{Type} \{t_\ell^i\}_{i\in I})}$$

It is easy to prove, by induction on $k$, that if $T$ is a finite subset of $\mathbb{T}_\kappa^\varsigma$, then $\sqcup_\kappa^\varsigma T$ (resp. $\sqcap_\kappa^\varsigma T$) is a least upper bound (resp. a greatest lower bound) for $T$ with respect to ${}_k{\leq}_\kappa^\varsigma$. As a result, it is also one with respect to $\leq_\kappa^\varsigma$. $\square$

In the following, we write $\leq$ instead of $\leq_\kappa^\varsigma$ when $\varsigma$ and $\kappa$ are irrelevant or can be inferred from the context. If $t \in \mathbb{T}_\kappa^{Type}$ is of the form $s(\bar{t}_p)$, then we let $t.p$ denote $t_p$. Similarly, if $t \in \mathbb{T}_\kappa^{Row}$ is of the form $(\bar{t}_\ell)$, then we let $t.\ell$ denote $t_\ell$. If $s \in \mathcal{S}_\kappa$ and $t \in \mathbb{T}_\kappa^{Type}$ is of the form $s'(\bar{t}_p)$, then we let $s \leq t$ stand for $s \leq_\kappa s'$.

## 2.3 Syntax of Terms and Constraints

For every sort $\varsigma$ and every kind $\kappa$, let $\mathcal{V}_\kappa^\varsigma$ be a distinct denumerable set of *variables*. Let $\alpha$ range over variables. Let a *filter* $L$ be a finite or cofinite subset of $\mathcal{L}$. Filters are machine representable and are preserved by finitary union, intersection, and complement. The syntax of *terms* and *constraints* is as follows:

$$\tau ::= \alpha \mid s(\bar{\tau}_p) \mid \partial\tau$$
$$C ::= \exists\alpha.C \mid C \wedge C \mid \textbf{true} \mid \textbf{false}$$
$$\mid \tau \leq \tau$$
$$\mid L : \tau \leq \tau$$
$$\mid L : s \leq \tau ? \tau \leq \tau$$

Terms are built out of variables, symbols, and the constructor $\partial$, which defines constant rows [13]. Constraint forms include existential quantification, conjunction, and *elementary* constraints, namely truth, falsity, subtyping constraints, and conditional subtyping constraints.

Figure 1 (resp. 2) defines the judgements $\vdash \tau : \varsigma$ (resp. $\vdash \tau : \kappa$), which means that the term $\tau$ has sort $\varsigma$ (resp. kind $\kappa$), and $\vdash^{\mathrm{OK}} C$ (resp. $\vdash_{\mathrm{OK}} C$), which means that the constraint $C$ is well-sorted (resp. well-kinded). We restrict our attention to well-sorted and well-kinded terms and constraints.

$$
\frac{\alpha \in \mathcal{V}_\kappa^\varsigma}{\vdash \alpha : \varsigma}
\qquad
\frac{\forall p \in a(s) \quad \vdash \tau_p : \mathrm{sort}(p)}{\vdash s(\bar{\tau}_p) : \mathit{Type}}
\qquad
\frac{\vdash \tau : \mathit{Type}}{\vdash \partial \tau : \mathit{Row}}
$$

$$
\frac{\vdash^{\mathrm{OK}} C}{\vdash^{\mathrm{OK}} \exists \alpha . C}
\qquad
\frac{\vdash^{\mathrm{OK}} C_1 \quad \vdash^{\mathrm{OK}} C_2}{\vdash^{\mathrm{OK}} C_1 \wedge C_2}
\qquad
\begin{array}{c}\vdash^{\mathrm{OK}} \mathbf{true} \\ \vdash^{\mathrm{OK}} \mathbf{false}\end{array}
$$

$$
\frac{\vdash \tau_1, \tau_2 : \mathit{Type}}{\vdash^{\mathrm{OK}} \tau_1 \le \tau_2}
\qquad
\frac{\vdash \tau_0, \tau_1, \tau_2 : \mathit{Row}}{\begin{array}{c}\vdash^{\mathrm{OK}} L : \tau_1 \le \tau_2 \\ \vdash^{\mathrm{OK}} L : s \le \tau_0 ? \tau_1 \le \tau_2\end{array}}
$$

**Figure 1. Well-sorted terms and constraints**

$$
\frac{\alpha \in \mathcal{V}_\kappa^\varsigma}{\vdash \alpha : \kappa}
\qquad
\frac{\begin{array}{c}s \in \mathcal{S}_\kappa\end{array}}{}
$$

$$
\frac{\alpha \in \mathcal{V}_\kappa^\varsigma}{\vdash \alpha : \kappa}
\qquad
\frac{\forall p \in a(s) \quad \vdash \tau_p : \mathrm{kind}(p)}{\vdash s(\bar{\tau}_p) : \kappa}
\qquad
\frac{\vdash \tau : \kappa}{\vdash \partial \tau : \kappa}
$$

$$
\frac{\vdash_{\mathrm{OK}} C}{\vdash_{\mathrm{OK}} \exists \alpha . C}
\qquad
\frac{\vdash_{\mathrm{OK}} C_1 \quad \vdash_{\mathrm{OK}} C_2}{\vdash_{\mathrm{OK}} C_1 \wedge C_2}
\qquad
\begin{array}{c}\vdash_{\mathrm{OK}} \mathbf{true} \\ \vdash_{\mathrm{OK}} \mathbf{false}\end{array}
$$

$$
\frac{\vdash \tau_1, \tau_2 : \kappa}{\vdash_{\mathrm{OK}} \tau_1 \le \tau_2 \atop \vdash_{\mathrm{OK}} L : \tau_1 \le \tau_2}
\qquad
\frac{s \text{ is a \textit{prime} element of } \mathcal{S}_\kappa \atop \vdash \tau_0, \tau_1, \tau_2 : \kappa}{\vdash_{\mathrm{OK}} L : s \le \tau_0 ? \tau_1 \le \tau_2}
$$

**Figure 2. Well-kinded terms and constraints**

If $\tau_1$ and $\tau_2$ are *types*, then the (standard) subtyping constraint $\tau_1 \le \tau_2$ may be understood as a requirement for $\tau_1$ to be a subtype of $\tau_2$. Constraints that involve *rows*, however, must be annotated with a filter. That is, if $\tau_1$ and $\tau_2$ are *rows*, then the constraint $L : \tau_1 \le \tau_2$ may be understood as a requirement for $\tau_1.\ell$ to be a subtype of $\tau_2.\ell$ *for every row label* $\ell \in L$. The introduction of filters compensates the omission of the standard row constructor ($\ell : \cdot ; \cdot$) by offering a way of *not* treating all row labels uniformly.

Our conditional constraints are exactly those studied in [11], extended with filters. Conditions must be of the form $s \le \tau$; this restriction allows our solver to delay inspecting a conditional constraint's conclusion until its condition must be satisfied, and is key to the algorithm's polynomial time complexity.

A symbol $s \in \mathcal{S}_\kappa$ is *prime* if and only if, for any finite subset $S$ of $\mathcal{S}_\kappa$, $s \le_\kappa (\sqcup_\kappa S)$ implies $\exists s' \in S \quad s \le_\kappa s'$. The requirement that $s$ be prime in Figure 2 slightly simplifies the constraint solver and leads to more aggressive constraint simplification algorithms, not discussed here; see [11]. It could be relaxed if required by some application.

## 2.4 Logical Interpretation

Terms and constraints are interpreted within the model under a *ground assignment* that gives meaning to their free variables.

**Definition 5** *A* ground assignment $\phi$ *is a total sort- and kind-preserving mapping from the variables into the model, that is, a total mapping from every* $\mathcal{V}_\kappa^\varsigma$ *into* $\mathbb{T}_\kappa^\varsigma$.

**Definition 6** *Ground assignments are extended to terms by*

$$
\phi(s(\bar{\tau}_p)) = s(\overline{\phi(\tau_p)})
$$
$$
\phi(\partial \tau) = (\overline{\phi(\tau)})
$$

The first clause of this definition is standard. The second clause interprets $\partial \tau$ under $\phi$ as the row that maps every row label to $\phi(\tau)$. If $\vdash \tau : \varsigma$ and $\vdash \tau : \kappa$ hold, then $\phi(\tau)$ is a member of $\mathbb{T}_\kappa^\varsigma$.

$$
\frac{\phi[\alpha \mapsto t] \vdash C}{\phi \vdash \exists \alpha . C}
\qquad
\frac{\phi \vdash C_1 \quad \phi \vdash C_2}{\phi \vdash C_1 \wedge C_2}
\qquad
\phi \vdash \mathbf{true}
$$

$$
\frac{\phi(\tau_1) \le \phi(\tau_2)}{\phi \vdash \tau_1 \le \tau_2}
\qquad
\frac{\forall \ell \in L \quad \phi(\tau_1).\ell \le \phi(\tau_2).\ell}{\phi \vdash L : \tau_1 \le \tau_2}
$$

$$
\frac{\forall \ell \in L \quad s \le \phi(\tau_0).\ell \Rightarrow \phi(\tau_1).\ell \le \phi(\tau_2).\ell}{\phi \vdash L : s \le \tau_0 ? \tau_1 \le \tau_2}
$$

**Figure 3. Constraint satisfaction**

**Definition 7** *Figure 3 defines the judgement* $\phi \vdash C$ *(read:* $\phi$ satisfies $C$*). The assertion* $C_1 \Vdash C_2$ *(read:* $C_1$ entails $C_2$*) holds if and only if, for every ground assignment* $\phi$, $\phi \vdash C_1$ *implies* $\phi \vdash C_2$. *Two constraints are* logically equivalent *if and only if they entail each other.*

4

The interpretation of constraints is straightforward. A constraint on rows is interpreted point-wise within its filter. A conditional constraint is interpreted as an implication.

The reader may notice that the constraints $\tau_1 \leq \tau_2$ and $L : \partial\tau_1 \leq \partial\tau_2$ are logically equivalent when $L$ is nonempty. So, it would be possible to suppress subtyping constraints on types altogether and to encode them as subtyping constraints on rows. It is convenient, however, to keep both; our constraint solver simplifies the latter into the former. In the case of conditional constraints, only constraints on rows were kept. These choices are somewhat arbitrary; other presentations would be possible.

## 3 Checking Satisfiability of Constraints

### 3.1 Representing Constraints

Thanks to $\alpha$-conversion of existentially bound variables and to scope extrusion, which allows rewriting $(\exists\alpha.C_1) \wedge C_2$ to $\exists\alpha.(C_1 \wedge C_2)$ when $\alpha$ does not appear free within $C_2$, every constraint can be put in prenex form. Given that $\exists\bar{\alpha}.C$ is satisfiable if and only if $C$ is, it is in fact sufficient to develop a satisfiability test for constraints that do not involve existential quantifiers.

We consider the conjunction operator $\wedge$ commutative, associative and idempotent, which allows viewing every constraint as a set of elementary constraints. This is not quite satisfactory, however, because such a representation could be very redundant. Indeed, given fixed rows $\tau_1$ and $\tau_2$, a constraint set might contain several—in fact, exponentially many—elementary constraints of the form $L : \tau_1 \leq \tau_2$. To remedy this problem, we identify constraints up to the following *fusion* laws:

$$(L_1 : \tau_1 \leq \tau_2) \wedge (L_2 : \tau_1 \leq \tau_2)$$
$$\equiv (L_1 \cup L_2) : \tau_1 \leq \tau_2$$
$$(L_1 : s \leq \tau_0 \,?\, \tau_1 \leq \tau_2) \wedge (L_2 : s \leq \tau_0 \,?\, \tau_1 \leq \tau_2)$$
$$\equiv (L_1 \cup L_2) : s \leq \tau_0 \,?\, \tau_1 \leq \tau_2$$

It is clear from Figure 3 that these laws preserve the meaning of constraints, which thus remains well-defined. Under these laws, a constraint $C$ may be viewed as a mapping from *edges* to *weights*, where an edge $e$ is

- one of **false** or $\tau_1 \leq \tau_2$, where $\vdash \tau_1, \tau_2 : \textit{Type}$; then, its weight $C(e)$ must be 0 or 1; or

- one of $\tau_1 \leq \tau_2$ or $s \leq \tau_0 \,?\, \tau_1 \leq \tau_2$, where $\vdash \tau_0, \tau_1, \tau_2 : \textit{Row}$; then, its weight $C(e)$ must be a filter $L$.

We write $C_1 \subseteq C_2$ if and only if some representatives of the constraints $C_1$ and $C_2$, viewed as sets of elementary constraints, are within the subset relation. Equivalently, viewing constraints as mappings from edges to weights, $\subseteq$ is the

pointwise ordering on mappings, provided weights are ordered by $0 \leq 1$ and set-theoretic inclusion of filters. If $c$ is an elementary constraint, the notation $c \in C$ stands for $\{c\} \subseteq C$; in particular, if $c$ is of the form, say, $L : \tau_1 \leq \tau_2$, then $c \in C$ is equivalent to $L \subseteq C(\tau_1 \leq \tau_2)$, that is, the edge $\tau_1 \leq \tau_2$ has weight $L$ *at least* in $C$.

### 3.2 Closing Constraints

We now define a procedure for determining whether a constraint is satisfiable. Following standard practice, it is presented as a closure computation, defined by the rules in Figure 4.

The first three rules deal with subtyping constraints between types; they implement transitive closure, structural decomposition, and failure, respectively. These rules are standard [6, 12].

The next three rules deal with subtyping constraints between rows. (TRANS-ROW) extends the usual transitivity rule by keeping track of filters: quite naturally, if $\tau_1 \leq \alpha$ holds at every row label $\ell \in L_1$ and if $\alpha \leq \tau_2$ holds at every row label $\ell \in L_2$, then transitivity applies at every row label $\ell \in L_1 \cap L_2$. (Recall that constraints between rows are interpreted pointwise.) (PROP-TYPE-ROW) extends (PROP) to the case where the sub-terms at $p$ of the types being decomposed are rows. The full filter $\mathcal{L}$ is used, because these rows should be comparable at every row label. (PROP-ROW-TYPE) decomposes a constraint between two uniform rows. Through these two rules, constraints between types may give rise to constraints between rows, and vice versa. These three rules are a key novelty of this paper. In particular, (TRANS-ROW) is a very natural generalization of (TRANS). None of the rules requires fresh variables to be allocated, contrary to the standard rules for row unification (see e.g. [14, appendix A]).

The last three rules deal with conditional constraints. (TRANS-CD-ROW) and (TRANS-CD-TYPE) are transitivity rules. They cause the lower bounds of the variable on which the condition bears (namely $\alpha$) to be examined. (FIRE) releases the conclusion into the pool as soon as it becomes evident that the condition is satisfied. These rules are a generalization of the standard closure rules for conditional constraints [11, appendix A] with filters.

Formally speaking, the closure rules define a monotonic function from constraints to constraints. This is made precise by the following definition.

**Definition 8** *Given a constraint $C$, let* $\text{close}(C)$ *be the least (w.r.t. $\subseteq$) constraint such that:*

- *$C \subseteq \text{close}(C)$; and*

- *if $C_1 \subseteq C$ and $C_1 \rightarrow C_2$ is an instance of one of the rules in Figure 4, then $C_2 \subseteq \text{close}(C)$.*

5

$$\begin{array}{lll}
(\textsc{Trans}) & \tau_1 \leq \alpha \wedge \alpha \leq \tau_2 & \rightarrow \quad \tau_1 \leq \tau_2 \\
(\textsc{Prop}) & s_1(\bar{\tau}_p^1) \leq s_2(\bar{\tau}_p^2) & \rightarrow \quad \tau_p^1 \stackrel{p}{\leq} \tau_p^2 \\
& & \quad \text{if } s_1 \leq s_2 \text{ and } p \in a(s_1) \cap a(s_2) \text{ and sort}(p) = \textit{Type} \\
(\textsc{Fail}) & s_1(\bar{\tau}_p^1) \leq s_2(\bar{\tau}_p^2) & \rightarrow \quad \textbf{false} \\
& & \quad \text{if } s_1 \not\leq s_2 \\
\\
(\textsc{Trans-Row}) & L_1 : \tau_1 \leq \alpha \wedge L_2 : \alpha \leq \tau_2 & \rightarrow \quad L_1 \cap L_2 : \tau_1 \leq \tau_2 \\
(\textsc{Prop-Type-Row}) & s_1(\bar{\tau}_p^1) \leq s_2(\bar{\tau}_p^2) & \rightarrow \quad \mathcal{L} : \tau_p^1 \stackrel{p}{\leq} \tau_p^2 \\
& & \quad \text{if } s_1 \leq s_2 \text{ and } p \in a(s_1) \cap a(s_2) \text{ and sort}(p) = \textit{Row} \\
(\textsc{Prop-Row-Type}) & L : \partial \tau_1 \leq \partial \tau_2 & \rightarrow \quad \tau_1 \leq \tau_2 \\
& & \quad \text{if } L \neq \varnothing \\
\\
(\textsc{Trans-Cd-Row}) & L_1 : s \leq \alpha \,?\, \tau_1 \leq \tau_2 \wedge L_2 : \partial \tau \leq \alpha & \rightarrow \quad L_1 \cap L_2 : s \leq \partial \tau \,?\, \tau_1 \leq \tau_2 \\
(\textsc{Trans-Cd-Type}) & L : s \leq \partial \alpha \,?\, \tau_1 \leq \tau_2 \wedge \tau \leq \alpha & \rightarrow \quad L : s \leq \partial \tau \,?\, \tau_1 \leq \tau_2 \\
(\textsc{Fire}) & L : s_1 \leq \partial s_2(\bar{\tau}_p) \,?\, \tau_1 \leq \tau_2 & \rightarrow \quad L : \tau_1 \leq \tau_2 \\
& & \quad \text{if } s_1 \leq s_2
\end{array}$$

**Figure 4. Constraint resolution**

If $C$ is well-sorted and well-kinded, then so is $\mathrm{close}(C)$. Furthermore, the closure rules are compatible with fusion equivalence: if $C_1 \equiv C_2$, then $\mathrm{close}(C_1) \equiv \mathrm{close}(C_2)$. This allows us to keep reasoning up to fusion equivalence.

Let us say that a row label $\ell$ is *apparent* in a set $L$ iff either $L$ is finite and $\ell \in L$, or $L$ is cofinite and $\ell \notin L$. Let us say that $\ell$ is apparent in a constraint $C$ iff $\ell$ is apparent in some $L$ that appears in $C$. Because the row labels apparent in $L_1 \cup L_2$ are those apparent in $L_1$ or $L_2$, this definition is compatible with fusion equivalence. The number of row labels apparent in a constraint $C$ is finite.

**Theorem 2** *Given a constraint $C$, the function* close *has a least fixed point containing $C$, called the* closure *of $C$.*

*Proof.* None of the rules in Figure 4 introduces new variables or symbols. The row labels apparent in $L_1 \cap L_2$ are those apparent in $L_1$ or $L_2$, so no rule makes new row labels apparent. Lastly, no rule builds new terms, except (TRANS-CD-TYPE), which builds the term $\partial \tau$ out of the term $\tau$. Because $\tau$ is of sort *Type*, while $\partial \tau$ is of sort *Row*, this process cannot be iterated: the number of terms in existence remains bounded. Given finite sets of symbols, variables, terms, and apparent row labels, only a finite number of elementary constraints may be built. Thus, the sequence $C, \mathrm{close}(C), \mathrm{close}^2(C), \ldots$ must converge. Its limit is the least fixed point of "close" containing $C$. $\square$

The following theorem implies that a constraint and its closure are logically equivalent.

**Theorem 3** *$C$ and* close$(C)$ *are logically equivalent.*

*Proof.* We have $\mathrm{close}(C) \supseteq C$, so $\mathrm{close}(C) \Vdash C$ holds. To establish the converse, we must check that, if $C_1 \rightarrow C_2$ is an instance of a closure rule, then $C_1 \Vdash C_2$ holds, that is, $\phi \vdash C_1$ implies $\phi \vdash C_2$.

◦ *Case* (TRANS) follows by transitivity of subtyping (Theorem 1).

◦ *Cases* (PROP), (FAIL) and (PROP-TYPE-ROW) follow from Definitions 4 and 6.

◦ *Case* (TRANS-ROW). Assume $\phi \vdash L_1 : \tau_1 \leq \alpha \wedge L_2 : \alpha \leq \tau_2$. Let $\ell \in L_1 \cap L_2$ (1). Because $\ell$ is a member of $L_1$, we have $\phi(\tau_1).r \leq \phi(\alpha).r$. Because $\ell$ is a member of $L_2$, we have $\phi(\alpha).r \leq \phi(\tau_2).r$. By transitivity, we obtain $\phi(\tau_1).r \leq \phi(\tau_2).r$. By discharging (1), we find $\phi \vdash L_1 \cap L_2 : \tau_1 \leq \tau_2$.

◦ *Case* (PROP-ROW-TYPE). Because $\phi(\partial \tau_i).\ell$ is $\phi(\tau_i)$, the statement $\phi \vdash L : \partial \tau_1 \leq \partial \tau_2$ is equivalent to $\forall \ell \in L \quad \phi(\tau_1) \leq \phi(\tau_2)$. Because $L$ is nonempty, the quantification is redundant; the statement implies $\phi \vdash \tau_1 \leq \tau_2$.

◦ *Case* (TRANS-CD-ROW). Assume $\phi \vdash L_1 : s \leq \alpha \,?\, \tau_1 \leq \tau_2 \wedge L_2 : \partial \tau \leq \alpha$. Let $\ell \in L_1 \cap L_2$ (1). Because $\ell$ is a member of $L_1$, we have $s \leq \phi(\alpha).\ell \Rightarrow \phi(\tau_1).\ell \leq \phi(\tau_2).\ell$ (2). Because $\ell$ is a member of $L_2$, we have $\phi(\tau) \leq \phi(\alpha).\ell$ (3). Assume $s \leq \phi(\partial \tau).\ell$, that is, $s \leq \phi(\tau)$ (4). It is straightforward to check that (4) and (3) imply $s \leq \phi(\alpha).\ell$, which, combined with (2), implies $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$. By successively discharging (4) and (1), we find $\phi \vdash L_1 \cap L_2 : s \leq \partial \tau \,?\, \tau_1 \leq \tau_2$.

◦ *Case* (TRANS-CD-TYPE). Assume $\phi \vdash L : s \leq \partial \alpha \,?\, \tau_1 \leq \tau_2 \wedge \tau \leq \alpha$. Let $\ell \in L$ (1). We have $s \leq \phi(\alpha) \Rightarrow \phi(\tau_1).\ell \leq \phi(\tau_2).\ell$ (2) and $\phi(\tau) \leq \phi(\alpha)$ (3).

Assume $s \leq \phi(\partial\tau).\ell$, that is, $s \leq \phi(\tau)$ (4). As above, (4) and (3) imply $s \leq \phi(\alpha)$, which, combined with (2), implies $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$. By successively discharging (4) and (1), we find $\phi \vdash L : s \leq \partial\tau ? \tau_1 \leq \tau_2$.

○ *Case* (FIRE). Assume $\phi \vdash L : s_1 \leq \partial s_2(\bar\tau_p) ? \tau_1 \leq \tau_2$. Let $\ell \in L$ (1). We have $s_1 \leq s_2(\phi(\bar\tau_p)) \Rightarrow \phi(\tau_1).\ell \leq \phi(\tau_2).\ell$. Because the rule carries the side-condition $s_1 \leq s_2$, this can be simplified to $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$. By discharging (1), we find $\phi \vdash L : \tau_1 \leq \tau_2$. □

A constraint $C$ is *closed* if and only if $C \supseteq \mathrm{close}(C)$. The following theorem shows that a closed constraint is satisfiable if and only if it does not contain **false**.

**Theorem 4** *If $C$ is closed and does not contain **false**, then $C$ is satisfiable.*

*Proof.* Let us first introduce some terminology. A ground assignment $\phi$ *$k$-satisfies* the constraint $\tau_1 \leq \tau_2$ if and only if $\phi(\tau_1) \ _k\!\leq \phi(\tau_2)$ holds. $\phi$ $k$-satisfies the constraint $L : \tau_1 \leq \tau_2$ if and only if $\phi(\tau_1).\ell \ _k\!\leq \phi(\tau_2).\ell$ holds for every $\ell \in L$. $\phi$ is a *$k$-candidate* for $C$ if and only if it $k$-satisfies every (non-conditional) subtyping constraint in $C$.

Let us continue with a couple of auxiliary definitions. When $\vdash \alpha : \textit{Type}$, let $\mathrm{lb}(\alpha)$ be the set of $\alpha$'s non-variable lower bounds in $C$, that is,

$$\mathrm{lb}(\alpha) = \{\tau ; \tau \leq \alpha \in C \wedge \tau \notin \mathcal{V}\}$$

When $\vdash \alpha : \textit{Row}$ and $\ell \in \mathcal{L}$, let $\mathrm{lb}_\ell(\alpha)$ be the set of $\alpha$'s lower bounds at $\ell$ in $C$, that is,

$$\mathrm{lb}_\ell(\alpha) = \{\tau ; \{\ell\} : \partial\tau \leq \alpha \in C\}$$

Recall that we reason up to fusion equivalence, so the condition $\{\ell\} : \partial\tau \leq \alpha \in C$ may be read $\ell \in C(\partial\tau \leq \alpha)$. If $\alpha$ has kind $\kappa$, then $\mathrm{lb}(\alpha)$ or $\mathrm{lb}_\ell(\alpha)$ (whichever is defined) is a finite set of terms of sort *Type* and of kind $\kappa$.

Now, consider the following system of equations, whose unknown is a ground assignment $\phi$:

$$\begin{aligned} \phi(\alpha) &= \ \sqcup\phi(\mathrm{lb}(\alpha)) & \text{if } \vdash \alpha : \textit{Type} \\ \phi(\alpha) &= \ \overline{(\sqcup\phi(\mathrm{lb}_\ell(\alpha)))} & \text{if } \vdash \alpha : \textit{Row} \end{aligned}$$

When $\alpha$ is a type variable, we require $\phi(\alpha)$ to be the least upper bound of the images through $\phi$ of $\alpha$'s lower bounds; when $\alpha$ is a row variable, we require $\phi(\alpha).\ell$ to be the least upper bound of the images through $\phi$ of $\alpha$'s lower bounds at $\ell$. The existence of least upper bounds is guaranteed by Theorem 1. Because $\mathrm{lb}(\alpha)$ contains terms of depth at least 1, the first equation is contractive; as a result, these equations admit a solution $\phi$. (See e.g. [9, Chapter 2] for a more detailed development, based on a metric argument.)

Let us now show that, for all $k \geq 0$, $\phi$ is a $k$-candidate for $C$. For $k = 0$, the result is immediate, since $ _0\!\leq$ is

everywhere true. Let us then assume that $\phi$ is a $k$-candidate and establish that it is in fact a $k + 1$-candidate. We pick an arbitrary (non-conditional) subtyping constraint in $C$ and proceed by cases according to its form. Because we reason up to fusion equivalence, we may restrict our attention to singleton filters.

○ *Case* $\alpha_1 \leq \alpha_2$. Because $C$ is closed with respect to (TRANS), we have $\mathrm{lb}(\alpha_1) \subseteq \mathrm{lb}(\alpha_2)$, which, by definition of $\phi$, implies $\phi(\alpha_1) \leq \phi(\alpha_2)$.

○ *Case* $\tau_1 \leq \alpha$, where $\tau_1 \notin \mathcal{V}$. Then, $\tau_1$ is a member of $\mathrm{lb}(\alpha)$. By definition of $\phi$, this implies $\phi(\tau_1) \leq \phi(\alpha)$.

○ *Case* $\tau_1 \leq \tau_2$, where $\tau_i$ is $s_i(\bar\tau_p^i)$ for each $i \in \{1, 2\}$. Because $C$ is closed with respect to (FAIL) and does not contain **false**, we have $s_1 \leq s_2$. Because $C$ is closed with respect to (PROP) and (PROP-TYPE-ROW) and because $\phi$ is a $k$-candidate, we have $\phi(\tau_p^1) \ _k^p\!\leq \phi(\tau_p^2)$ for every $p \in a(s_1) \cap a(s_2)$. As a result, we have $\phi(\tau_1) \ _{k+1}\!\leq \phi(\tau_2)$.

○ *Case* $\alpha \leq \tau_2$, where $\tau_2 \notin \mathcal{V}$. Assume $\tau_1 \in \mathrm{lb}(\alpha)$. Because $C$ is closed with respect to (TRANS), we have $\tau_1 \leq \tau_2 \in C$. According to the previous case, this implies $\phi(\tau_1) \ _{k+1}\!\leq \phi(\tau_2)$. Because this holds for an arbitrary choice of $\tau_1$ within $\mathrm{lb}(\alpha)$, we have $\sqcup\phi(\mathrm{lb}(\alpha)) \ _{k+1}\!\leq \phi(\tau_2)$, that is, $\phi(\alpha) \ _{k+1}\!\leq \phi(\tau_2)$.

○ *Case* $\{\ell\} : \alpha_1 \leq \alpha_2$. Because $C$ is closed with respect to (TRANS-ROW), we have $\mathrm{lb}_\ell(\alpha_1) \subseteq \mathrm{lb}_\ell(\alpha_2)$, which, by definition of $\phi$, implies $\phi(\alpha_1).\ell \leq \phi(\alpha_2).\ell$.

○ *Case* $\{\ell\} : \partial\tau_1 \leq \alpha$. Then, $\tau_1$ is a member of $\mathrm{lb}_\ell(\alpha)$. By definition of $\phi$, this implies $\phi(\tau_1) \leq \phi(\alpha).\ell$, that is, $\phi(\partial\tau_1).\ell \leq \phi(\alpha).\ell$.

○ *Case* $\{\ell\} : \partial\tau_1 \leq \partial\tau_2$. Because $C$ is closed with respect to (PROP-ROW-TYPE), we have $\tau_1 \leq \tau_2 \in C$. According to the four previous cases, this implies $\phi(\tau_1) \ _{k+1}\!\leq \phi(\tau_2)$, that is, $\phi(\partial\tau_1).\ell \ _{k+1}\!\leq \phi(\partial\tau_2).\ell$.

○ *Case* $\{\ell\} : \alpha \leq \partial\tau_2$. Assume $\tau_1 \in \mathrm{lb}_\ell(\alpha)$. Then, we have $\{\ell\} : \partial\tau_1 \leq \alpha \in C$. Because $C$ is closed with respect to (TRANS-ROW), $\{\ell\} : \partial\tau_1 \leq \partial\tau_2 \in C$ must hold. According to the previous case, this implies $\phi(\tau_1) \ _{k+1}\!\leq \phi(\partial\tau_2).\ell$. Because this holds for an arbitrary choice of $\tau_1$ within $\mathrm{lb}_\ell(\alpha)$, we have $\sqcup\phi(\mathrm{lb}_\ell(\alpha)) \ _{k+1}\!\leq \phi(\partial\tau_2).\ell$, that is, $\phi(\alpha).\ell \ _{k+1}\!\leq \phi(\partial\tau_2).\ell$.

Our case analysis, and our induction on $k$, are now complete; we have established that $\phi$ satisfies every (non-conditional) subtyping constraint in $C$. Let us now prove that $\phi$ also satisfies the conditional constraints in $C$. We again pick such a constraint and proceed by cases according to its form, assuming a singleton filter.

○ *Case* $\{\ell\} : s_1 \leq \partial s_2(\bar\tau_p) ? \tau_1 \leq \tau_2$. Assume $s_1 \leq \phi(\partial s_2(\bar\tau_p)).\ell$. This may be written $s_1 \leq \phi(s_2(\bar\tau_p))$, that is, $s_1 \leq s_2$. Then, because $C$ is closed with respect to (FIRE), we have $\{\ell\} : \tau_1 \leq \tau_2 \in C$. Because $\phi$ satisfies every (non-conditional) subtyping constraint in $C$, this im-

plies $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$.

○ *Case* $\{\ell\} : s \leq \partial\alpha ? \tau_1 \leq \tau_2$. Assume $s \leq \phi(\partial\alpha).\ell$, that is, $s \leq \phi(\alpha)$. By definition of $\phi$, this is $s \leq \sqcup \phi(\mathrm{lb}(\alpha))$. Because the head symbol of the least upper bound of a set of ground terms is the least upper bound of their head symbols, and because $s$ must be a prime element of its symbol lattice, this implies $s \leq \phi(\tau)$ for some $\tau \in \mathrm{lb}(\alpha)$. Thus, we have $\tau \leq \alpha \in C$. Because $C$ is closed with respect to (TRANS-CD-TYPE), we have $\{\ell\} : s \leq \partial\tau ? \tau_1 \leq \tau_2 \in C$. According to the previous case, $\phi$ satisfies this constraint. Given $s \leq \phi(\tau)$, this implies $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$.

○ *Case* $\{\ell\} : s \leq \alpha ? \tau_1 \leq \tau_2$. Assume $s \leq \phi(\alpha).\ell$. By definition of $\phi$, this is $s \leq \sqcup \phi(\mathrm{lb}_\ell(\alpha))$. As above, this implies $s \leq \phi(\tau)$ for some $\tau \in \mathrm{lb}_\ell(\alpha)$. Thus, we have $\{\ell\} : \partial\tau \leq \alpha \in C$. Because $C$ is closed with respect to (TRANS-CD-ROW), we have $\{\ell\} : s \leq \partial\tau ? \tau_1 \leq \tau_2 \in C$. According to the previous cases, $\phi$ satisfies this constraint. Given $s \leq \phi(\tau)$, this implies $\phi(\tau_1).\ell \leq \phi(\tau_2).\ell$. □

Together, these results show that constraint satisfiability is decidable: a constraint is satisfiable if and only if its closure (which may be computed in a finite number of steps) does not contain **false**.

Theorem 3 shows that a constraint and its closure are *logically* equivalent, whereas proving that they are *satisfaction* equivalent (i.e. one is satisfiable if and only if the other is satisfiable) would suffice to establish decidability. The latter approach is followed by Palsberg and Zhao [7, Lemma 5.4]. It is, however, much weaker, because it leads to a *non-compositional* type inference system. Indeed, constraint generation and constraint resolution can be interleaved only if the latter preserves the meaning of constraints, that is, the set of their solutions. If it doesn't, then resolution must be performed *after* generation is complete, that is, after the whole program has been analyzed. A compositional approach, where constraint generation may be freely interleaved with resolution and simplification, is much more desirable, for efficiency and modularity reasons.

Palsberg and Zhao's notion of "esat-closure" eagerly simplifies constraints of the form $V \oplus V' \leq [\ell : U]^\rightarrow$ into either $V \leq [\ell : U]^\rightarrow$ or $V' \leq [\ell : U]^\rightarrow$. When the context does not allow determining which of these two choices is right, an arbitrary decision is made. However, in such a situation, there is no principal choice—that is, each of the two possible choices rules out some solutions—so the solution set is not preserved. Our approach, on the other hand, is to *delay* this choice until it can be safely resolved. We achieve this via conditional constraints; see the type scheme ascribed to symmetric record concatenation in Section 4.

### 3.3 Description of the Algorithm

We now describe an algorithm that, given a constraint $C_0$, fails if it is unsatisfiable, and computes its closure oth-

erwise. The algorithm uses an unordered *queue $Q$* of pending elementary constraints. It also maintains a *current constraint $C$*, represented as an association table from edges to weights, as suggested in Section 3.1. $C$ is initially **true**, that is, every edge has weight 0 or $\varnothing$, as appropriate.

$$
\boxed{
\begin{aligned}
&\textbf{if } c \text{ is } e \textbf{ and } C(e) = 0 \\
&\textbf{then } \begin{cases} C(e) \leftarrow 1 \\ Q \leftarrow Q, c \end{cases} \\[1em]
&\textbf{else if } c \text{ is } L : e \\
&\textbf{then } \begin{cases} \textbf{let } L_1 = L \setminus C(e) \\ \textbf{if } L_1 \neq \varnothing \\ \textbf{then } \begin{cases} C(e) \leftarrow L \cup C(e) \\ Q \leftarrow Q, L_1 : e \end{cases} \end{cases}
\end{aligned}
}
$$

**Figure 5. Algorithm** INSERT($c$)

Pseudocode for the procedure INSERT($c$), which inserts an elementary constraint $c$ into $C$ and $Q$, is given in Figure 5. In "**if** $c$ is *pattern* **then** *command*", the meta-variables that occur within the pattern are bound in the command. If $c$ is a constraint on types, INSERT adds it to $C$ and $Q$, unless it has already been processed before. If, on the other hand, $c$ is a constraint on rows, of the form $L : e$, then it is first rewritten to $L_1 : e$, where $L_1$ is $L \setminus C(e)$, that is, $L$ minus the edge's current weight. The rewritten constraint is then added to $C$ and $Q$ if it is nontrivial. The point of rewriting the constraint in such a way is to avoid any duplication of effort; the complexity analysis (Section 3.5) relies on this optimization.

Using INSERT, we define the procedure PROCESS($c$), given in Figure 6, whose purpose is to enumerate and insert the immediate consequences of a newly discovered constraint $c$. The procedure explores all new manners in which a closure rule may become applicable due to the addition of $c$. It is mostly a paraphrase of Figure 4.

Figure 7 defines the main procedure. Starting with an empty current constraint $C$ and an empty queue $Q$, the function uses INSERT to schedule every element of $C_0$ for consideration. Then, as long as some constraint remains in the queue, it is processed, causing its consequences to be inserted into the queue, unless they were known already. The algorithm stops when INSERT fails or when the queue becomes empty.

### 3.4 Properties of the Algorithm

**Theorem 5 (Soundness)** *The main **while** loop has invariant property $C_0 \subseteq C \subseteq \mathrm{close}^\infty(C_0)$.*

*Proof.* (Sketch.) The initialization step establishes $C \equiv C_0$, so this property is initially satisfied. To show that

$$\textbf{if } c \textbf{ is } \tau_1 \leq \alpha \qquad\qquad (1)$$
$$\textbf{then for each } \tau_2 \textbf{ such that } C(\alpha \leq \tau_2) = 1$$
$$\quad \textbf{do } \text{INSERT}(\tau_1 \leq \tau_2)$$

$$\textbf{if } c \textbf{ is } \alpha \leq \tau_2 \qquad\qquad (2)$$
$$\textbf{then for each } \tau_1 \textbf{ such that } C(\tau_1 \leq \alpha) = 1$$
$$\quad \textbf{do } \text{INSERT}(\tau_1 \leq \tau_2)$$

$$\textbf{if } c \textbf{ is } L_1 : \tau_1 \leq \alpha \qquad\qquad (3)$$
$$\textbf{then for each } \tau_2$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L_2 = C(\alpha \leq \tau_2) \\ \text{INSERT}(L_1 \cap L_2 : \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } L_2 : \alpha \leq \tau_2 \qquad\qquad (4)$$
$$\textbf{then for each } \tau_1$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L_1 = C(\tau_1 \leq \alpha) \\ \text{INSERT}(L_1 \cap L_2 : \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } s_1(\bar{\tau}_p^1) \leq s_2(\bar{\tau}_p^2) \qquad\qquad (5)$$
$$\textbf{then } \begin{cases} \textbf{if } s_1 \not\leq s_2 \textbf{ then fail} \\ \textbf{else for each } p \in a(s_1) \cap a(s_2) \\ \quad \textbf{do } \begin{cases} \textbf{if } \text{sort}(p) = \textit{Type} \\ \textbf{then } \text{INSERT}(\tau_p^1 \overset{p}{\leq} \tau_p^2) \\ \textbf{else } \text{INSERT}(\mathcal{L} : \tau_p^1 \overset{p}{\leq} \tau_p^2) \end{cases} \end{cases}$$

$$\textbf{if } c \textbf{ is } L : \partial \tau_1 \leq \partial \tau_2 \qquad\qquad (6)$$
$$\textbf{then } \text{INSERT}(\tau_1 \leq \tau_2)$$

$$\textbf{if } c \textbf{ is } L_1 : s \leq \alpha \, ? \, \tau_1 \leq \tau_2 \qquad\qquad (7)$$
$$\textbf{then for each } \tau$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L_2 = C(\partial \tau \leq \alpha) \\ \text{INSERT}(L_1 \cap L_2 : s \leq \partial \tau \, ? \, \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } L_2 : \partial \tau \leq \alpha \qquad\qquad (8)$$
$$\textbf{then for each } s, \tau_1, \tau_2$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L_1 = C(s \leq \alpha \, ? \, \tau_1 \leq \tau_2) \\ \text{INSERT}(L_1 \cap L_2 : s \leq \partial \tau \, ? \, \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } L : s \leq \partial \alpha \, ? \, \tau_1 \leq \tau_2 \qquad\qquad (9)$$
$$\textbf{then for each } \tau$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L = C(\tau \leq \alpha) \\ \text{INSERT}(L : s \leq \partial \tau \, ? \, \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } \tau \leq \alpha \qquad\qquad (10)$$
$$\textbf{then for each } s, \tau_1, \tau_2$$
$$\quad \textbf{do } \begin{cases} \textbf{let } L = C(s \leq \partial \alpha \, ? \, \tau_1 \leq \tau_2) \\ \text{INSERT}(L : s \leq \partial \tau \, ? \, \tau_1 \leq \tau_2) \end{cases}$$

$$\textbf{if } c \textbf{ is } L : s_1 \leq \partial s_2(\bar{\tau}_p) \, ? \, \tau_1 \leq \tau_2 \textbf{ and } s_1 \leq s_2 \quad (11)$$
$$\textbf{then } \text{INSERT}(L : \tau_1 \leq \tau_2)$$

**Figure 6. Algorithm** $\text{PROCESS}(c)$

$$C \leftarrow \textbf{true}$$
$$Q \leftarrow \text{empty}$$
$$\textbf{for each } c \in C_0$$
$$\quad \textbf{do if } c \textbf{ is false then fail else } \text{INSERT}(c)$$
$$\textbf{while } Q \textbf{ is nonempty}$$
$$\quad \textbf{do } \begin{cases} \text{extract } c \text{ out of } Q \\ \text{PROCESS}(c) \end{cases}$$

**Figure 7. Algorithm** $\text{CLOSURE}(C_0)$

it is preserved throughout execution, one first establishes another invariant property, namely $Q \subseteq C$, by inspection of the definition of INSERT. This ensures that, when PROCESS($c$) is invoked, $c \in C$ holds. As a result, every new constraint inserted by PROCESS is a member of close($C$). The result follows. $\qquad\square$

**Theorem 6 (Completeness)** *The main **while** loop has invariant property* $\text{close}(C \setminus Q) \subseteq C$.

*Proof.* (Sketch.) The initialization step establishes $C \equiv Q$, so this property is initially satisfied. Let $(C_i, Q_i)$ be the algorithm's state after $i$ runs through the main **while** loop. Write $c.C$ for the constraints that can be derived by applying a closure rule to $c$ and (if it is a binary rule) to some element of $C$. Assume the invariant holds at step $i$. Then:

$$\begin{aligned} & \text{close}(C_{i+1} \setminus Q_{i+1}) \\ = \;& \text{close}(\{c\} \cup (C_i \setminus Q_i)) \\ & \text{because PROCESS does not affect } C \setminus Q \\ = \;& \{c\} \cup c.(C_i \setminus Q_i) \cup \text{close}(C_i \setminus Q_i) \\ \subseteq \;& c.C_i \cup C_i \\ & \text{because } c \in C_i \\ & \text{by monotonicity of } c.C \text{ w.r.t. } C \\ & \text{by the invariant at step } i \\ = \;& C_{i+1} \\ & \text{by definition of PROCESS} \end{aligned}$$

which establishes the invariant at step $i + 1$. $\qquad\square$

When the algorithm succeeds, $Q$ is empty, so, by Theorem 6, $\text{close}(C) \subseteq C$ holds—that is, $C$ is closed. Together with Theorem 5, this shows that $C$ is then the closure of the initial constraint $C_0$. Conversely, when the algorithm fails, $\text{close}(C)$ contains **false**, so, by Theorem 5, the closure of $C_0$ contains **false** as well. Thus, the algorithm succeeds if and only if the closure of $C_0$ does not contain **false**, that is—according to Theorems 3 and 4—if and only if $C_0$ is satisfiable.

### 3.5 Complexity Analysis

We assess the algorithm's time and space complexity in terms of two parameters, namely $n$, the size of $C_0$, and $m$,

the number of row labels apparent in $C_0$. Note that, for $n$ to be properly defined, we must view $C_0$ as a multiset of elementary constraints; we do *not* apply the fusion laws to it. The parameter $m$ is bounded by $n$, but is typically much smaller in practice; hence, it is worth distinguishing between them.

By examination of the closure rules (Figure 4), one finds that the number of edges whose weight may become non-null at some point is only $O(n^2)$. The crucial point is that even though new conditional constraints of the form $L : s \leq \tau_0 ? \tau_1 \leq \tau_2$ may be created, no new triples $(s, \tau_1, \tau_2)$ appear, so the number of distinct such triples remains $O(n)$. Given that the number of choices for $\tau_0$ is also $O(n)$, this leads us to a quadratic number of edges of the form $s \leq \tau_0 ? \tau_1 \leq \tau_2$. Subtyping edges are clearly at most quadratic in number as well.

The association table from edges to weights is implemented as an array, whose space usage is thus $O(n^2)$. Access and update operations have $O(1)$ time complexity. We conduct our analysis based on these assumptions, although it would be desirable, in practice, to use a more space-efficient data structure, possibly at the expense of an additional $O(\log n)$ time penalty. Insertion and removal in $Q$ are performed in constant time.

We assume that $\mathcal{L}$ is totally ordered, so filters may be represented as a pair of a Boolean flag and a balanced binary tree whose keys are row labels. A filter may be checked for emptiness in constant time. Furthermore, membership $\ell \in L$ and insertion $\{\ell\} \cup L$ have $O(\log m)$ time complexity. As a result, union $L_1 \cup L_2$, intersection $L_1 \cap L_2$ and difference $L_1 \setminus L_2$ may be computed in time $O(1 + |L_1| \log m)$, where $|L|$ stands for the cardinal of $L$ if $L$ is finite, and for that of its complement if it is cofinite.

If $c$ is a constraint on types, then $\text{INSERT}(c)$ runs in time $O(1)$. If, on the other hand, $c$ is a constraint on rows, of the form $L : e$, then it runs in time $O(1 + |L| \log m)$.

A constraint of the form $e$ can appear at most once in $Q$, because $\text{INSERT}$ checks whether $C(e)$ is $0$ before appending $e$ to $Q$. Furthermore, by definition of $\text{INSERT}$, if constraints of the form $L_1 : e, L_2 : e, \dots, L_k : e$ successively appear in the queue, then the filters $(L_i)_{1 \leq i \leq k}$ must be non-empty and pairwise disjoint. This implies, in particular, that at most one $L_i$ is cofinite. All others must be finite, and because they are pairwise disjoint, there are at most $m$ of them. As a result, the sum $\Sigma_{1 \leq i \leq k} O(1 + |L_i| \log m)$ is bounded by $O(m \log m)$. This remark will be used shortly.

Let us now examine block 1 in Figure 6, which implements half of (TRANS). The number of distinct constraints of the form $\tau_1 \leq \alpha$ is $O(n^2)$, so this block is executed at most $O(n^2)$ times. There are $O(n)$ distinct terms $\tau_2$. Lastly, looking up $C$ and invoking $\text{INSERT}$ takes time $O(1)$ in this case, so the total cost for this block is $O(n^3)$. Block 2, which is symmetric and implements the other half

of (TRANS), has the same cost. Indeed, $O(n^3)$ is the cost usually associated with transitive closure.

Let us now examine block 3, which implements half of (TRANS-ROW). What is the cost of its inner **do** statement? If $L_1$ is finite, then $|L_1 \cap L_2| \leq |L_1|$ holds, so the cost of the call to $\text{INSERT}$ is $O(1 + |L_1| \log m)$, and the cost of the **do** statement is $O(1 + |L_1| \log m)$ as well. If, on the other hand, $L_1$ is cofinite, then we may bound the statement's cost by $O(m \log m)$. According to the remark above, if we now let $L_1$ vary while $\tau_1$, $\tau_2$ and $\alpha$ remain fixed, the statement's cumulative cost is only $O(m \log m)$. Lastly, letting $\tau_1$, $\tau_2$ and $\alpha$ vary, we find that the total cost for block 3 is $O(n^3 m \log m)$. Block 4 is symmetric.

We now come to block 5, which implements (PROP) and (PROP-TYPE-ROW). Assuming that $\text{sort}(p)$ may be looked up in constant time, the cost of the inner **do** statement is $O(1)$. The cardinality of the set $a(s_1) \cap a(s_2)$ is bounded by a constant, because Definition 2 specifies that the set of symbols in the (fixed) ground signature is finite. Assuming that this set can be computed in constant time, the cost of the **for** loop is still $O(1)$. Let us assume that symbols can be compared in constant time. Because the number of distinct constraints of the form $s_1(\bar{\tau}_p^1) \leq s_2(\bar{\tau}_p^2)$ is $O(n^2)$, the total cost for block 5 is $O(n^2)$ as well.

Block 6, which implements (PROP-ROW-TYPE), has cost $O(mn^2)$. This again follows from the fact that, when $L$ varies, its successive values are, on the one hand, at most one cofinite filter, and on the other hand, a sequence of pairwise disjoint, finite filters; so $L$ may only assume $O(m)$ successive values.

The analysis for blocks 7, 8, 9 and 10, which implement (TRANS-CD-ROW) and (TRANS-CD-TYPE), is similar to that for block 3: each of these blocks has total cost $O(n^3 m \log m)$. The analysis for block 11 is analogous to that for block 6: it has total cost $O(n^2 m \log m)$. To obtain these bounds, one must recall that the number of triples $(s, \tau_1, \tau_2)$ in conditional constraints is bounded by $O(n)$.

So far, we have measured the cost of executing each block, when the **if** statement that governs it is taken. There remains to assess the cumulative cost of going through these **if** statements, which is non-null even when the branch is not taken. It is clear that running one constraint $c$ through all **if** statements costs $O(1)$ time. Because the number of constraints that can be queued is $O(mn^2)$, the cumulative cost of going through these statements is $O(mn^2)$ as well. To sum up, we have established

**Theorem 7** *The algorithm runs in time* $O(n^3 m \log m)$.

## 4 Typing Record Operations

In this section, we use the ground signature defined in Example 1. We instantiate the parametric constraint-based

10

type inference framework $\text{HM}(X)$ [5] with our constraint language. Then, we extend it with primitive operations on records, to which we assign *constrained type schemes*, as follows. Let $\alpha$ range over variables of sort *Type* and kind *type*. Let $\varphi$ range over variables of sort *Row* and kind *field*. The *empty record* has type $\{\partial\mathsf{Abs}\}$. *Access* $(.\ell)$ has type $\forall\alpha\varphi[\{\ell\} : \varphi \leq \partial(\mathsf{Pre}\,\alpha)].\{\varphi\} \to \alpha$. Notice how a singleton filter is used to extract information about field $\ell$ alone. *Non-strict extension* $(+\ell)$, which subsumes extension and update, has type

$$\forall\alpha\varphi_1\varphi_2[\{\ell\} : \partial(\mathsf{Pre}\,\alpha) \leq \varphi_2 \wedge (\mathcal{L} \setminus \{\ell\}) : \varphi_1 \leq \varphi_2].$$
$$\{\varphi_1\} \to \alpha \to \{\varphi_2\}$$

Here, the first constraint uses a singleton filter to indicate that the field $\ell$ is present with type $\alpha$ in the new record, while the second constraint uses a cosingleton filter to indicate that all fields other than $\ell$ have the same status as in the original record. *Symmetric concatenation* $(+)$ has type

$$\forall\varphi_1\varphi_2\varphi_3[\ \mathcal{L} : \mathsf{Abs} \leq \varphi_1 ? \varphi_2 \leq \varphi_3$$
$$\wedge \mathcal{L} : \mathsf{Abs} \leq \varphi_2 ? \varphi_1 \leq \varphi_3$$
$$\wedge \mathcal{L} : \mathsf{Pre} \leq \varphi_1 ? \varphi_2 \leq \mathsf{Abs}].$$
$$\{\varphi_1\} \to \{\varphi_2\} \to \{\varphi_3\}$$

The first two constraints encode the semantics of record concatenation: if a field is missing from one argument, it is read from the other. The third constraint prevents a field from being present in both arguments at once. All constraints carry the full filter $\mathcal{L}$, because concatenation behaves uniformly with respect to all field labels. Assigning an accurate type to *asymmetric* record concatenation requires a slight extension of our framework; see Section 5.3.

It is straightforward to give a semantics to these operations and to prove that these types are correct with respect to it. We omit this step by lack of space.

Our approach presents no novelty with respect to previous work [14, 11], except in the use of our new constraint language. Because the types above have bounded size, the size of a constraint associated with a monomorphic program is linear in the size of the program. The row labels apparent in the constraint are those apparent in the program text. Thus, we are able to infer types in time $O(n^3 m \log m)$, where $n$ is the program size and $m$ is the number of record labels in the program. The addition of **let**-polymorphism invalidates this result in principle but is known to have little practical impact on performance, provided effective constraint simplification algorithms are available.

## 5 Extensions

Sections 3.5 and 4 are based on the assumption that a filter is encoded as a pair of a Boolean flag and a finite set of row labels. However, the definition and proof of the constraint solving algorithm only require filters to be preserved by Boolean operations and to enjoy a decidable emptiness test. Sections 5.1 and 5.2 exploit this remark and suggest richer filter languages. Section 5.3 suggests another extension, based on a more liberal sorting discipline.

### 5.1 Hierarchizing Row Labels

Imagine $\mathcal{L}$ is *finite* and equipped with a partial order $\preceq$. Imagine the programmer specifies $(\mathcal{L}, \preceq)$, in a modular manner, by declaring new elements and edges at the beginning of every program module. Thus, $(\mathcal{L}, \preceq)$ forms a finite but *extensible* hierarchy. Let the *cone* generated by $\ell$ be the set of $\ell$'s lower bounds with respect to $\preceq$. Let *filter expressions* $\lambda$ consist of singletons, cones (represented by their generator), and Boolean combinations thereof:

$$\lambda ::= \{\ell\} \mid \downarrow\ell \mid \lambda \cup \lambda \mid \neg\lambda$$

Every filter expression $\lambda$ may be interpreted as a filter in the obvious manner. Replace filters with filter expressions in the definition of constraints: the meaning of constraints is now parameterized by $(\mathcal{L}, \preceq)$.

If $(\mathcal{L}, \preceq)$ was known entirely when solving constraints, there would be little point in this extension, because $\downarrow\ell$ would be syntactic sugar for a union of singletons. (It would be more concise, though, which may be interesting in practice.) However, the row label hierarchy can be specified modularly, which means that, when a program module $M$ is being examined in isolation, $(\mathcal{L}, \preceq)$ is only known to be *some* extension of the hierarchy $(\mathcal{L}_M, \preceq_M)$ declared in $M$. Thus, one must ensure that the constraints associated with $M$ are satisfiable under *every* extension of $(\mathcal{L}_M, \preceq_M)$. To this end, one need not modify our algorithm; it suffices to plug in a test for emptiness that determines whether a given filter expression is empty under every extension of the current hierarchy. For such a test to be decidable, one must likely restrict the way new elements and edges are declared; we leave this issue for future work.

How may such an extension be useful? Consider a type system that gathers information about every object's actual class in a Java-like programming language. Let $(\mathcal{L}, \preceq)$ reflect the class and interface hierarchy. Then, the **new** operator may be described using a singleton filter, because it creates objects of a known class. The **instanceof** operator may be described using a cone filter, because it selects objects whose actual class is *any* subclass of a known class. Using this mechanism, uncaught exception analyses based on rows (see e.g. [8]) could be extended to handle languages where exceptions form a hierarchy, such as Java, while preserving their ability to analyze program modules in isolation. Naturally, this extension is speculative; its details remain to be worked out.

## 5.2 Structuring Individuals

Let us say that a filter is a set of *individuals*. So far, we have considered individuals to be row labels, that is, atoms. What if instead individuals were structured entities?

For instance, let individuals be $k$-tuples of row labels, for some fixed positive integer $k$. Let filters be finite unions of Cartesian products of the form $\mathcal{L}^i \times L \times \mathcal{L}^{k-1-i}$, where $0 \leq i < k$ and $L$ is a finite or cofinite subset of $\mathcal{L}$, as before. Then, filters are preserved by intersection and complement and enjoy a decidable emptiness test. Without modifying the constraint solving procedure, we obtain a constraint language that offers $k$-*dimensional rows*, a concept previously studied by Rémy. Rémy suggested using a 2-dimensional row (that is, a 2-dimensional array of types), indexed on the one hand by method names and the other hand by class names, to keep track of the type of every method at every class in an object-oriented language. Using filters of the form $L \times \mathcal{L}$ or $\mathcal{L} \times L$, it is possible to update an entire line or an entire column of the array at once, which allows assigning types to the operations that create new methods or new classes.

As another instance, let individuals be *trees*; let filters be some class of *tree automata* that is preserved by Boolean operations and enjoys a decidable emptiness test. We obtain a new flavor of rows, which allows encoding certain functions from trees into types. This may open an avenue towards full type inference for programming languages dedicated to manipulating XML structures, such as XDuce, that have *regular expression types* [2]. Currently, XDuce only has some (local) type inference.

Again, an in-depth exploration of these extensions must be left for future research. Still, we believe it is pleasing and interesting that the algorithm should exhibit such generality.

## 5.3 More Row Terms

Figure 1 requires types of the form $s(\bar{\tau}_p)$ to have sort *Type*, which means that row terms must be variables or of the form $\partial \tau$. This is indeed sufficient for many applications. However, it is sometimes desirable to give meaning to row terms of the form $s(\bar{\tau}_p)$. Rémy [13] pointed out that every type constructor $s$ of sort $Type^k \Rightarrow Type$ may be viewed, at the same time, as a row constructor of sort $Row^k \Rightarrow Row$. In the logical model, the effect of $s$ on rows is then defined as the point-wise extension of its effect on types. This extension yields extra expressiveness, which is useful to infer types for language features such as *asymmetric* record concatenation and so-called *first-class messages* [11]. We anticipate no difficulty in adapting our results to this more liberal sorting discipline. The algorithm must then be extended with more closure rules; we have avoided this complication for the sake of clarity.

## References

[1] Luca Cardelli and John Mitchell. Operations on records. In Carl A. Gunter and John C. Mitchell, editors, *Theoretical Aspects of Object-Oriented Programming: Types, Semantics, and Language Design*. MIT Press, 1994. URL: `http://research.microsoft.com/Users/luca/Papers/Records.ps`.

[2] Haruo Hosoya and Benjamin C. Pierce. Regular expression pattern matching for XML. *Journal of Functional Programming*, 2002. To appear. URL: `http://www.kurims.kyoto-u.ac.jp/~hahosoya/papers/tapat-full.ps`.

[3] Dexter Kozen, Jens Palsberg, and Michael I. Schwartzbach. Efficient recursive subtyping. *Mathematical Structures in Computer Science*, 5(1):113–125, 1995. URL: `http://www.cs.purdue.edu/homes/palsberg/paper/mscs95-kps.ps.gz`.

[4] David McAllester. On the complexity analysis of static analyses. *Journal of the ACM*, 49(4):512–537, July 2002. URL: `http://www.autoreason.com/sas01.ps`.

[5] Martin Odersky, Martin Sulzmann, and Martin Wehr. Type inference with constrained types. *Theory and Practice of Object Systems*, 5(1):35–55, 1999. URL: `http://www.comp.nus.edu.sg/~sulzmann/publications/tapos.ps`.

[6] Jens Palsberg, Mitchell Wand, and Patrick M. O'Keefe. Type inference with non-structural subtyping. *Formal Aspects of Computing*, 9:49–67, 1997. URL: `http://www.cs.purdue.edu/homes/palsberg/paper/fac97.ps.gz`.

[7] Jens Palsberg and Tian Zhao. Efficient type inference for record concatenation and subtyping. In *IEEE Symposium on Logic in Computer Science (LICS)*, pages 125–136, July 2002. URL: `http://www.cs.purdue.edu/homes/palsberg/paper/lics02.ps.gz`.

[8] François Pessaux and Xavier Leroy. Type-based analysis of uncaught exceptions. *ACM Transactions on Programming Languages and Systems*, 22(2):340–377, 2000. URL: `http://pauillac.inria.fr/~xleroy/publi/exceptions-toplas.ps.gz`.

[9] François Pottier. Type inference in the presence of subtyping: from theory to practice. Research Report 3483, INRIA, September 1998. URL: `ftp://ftp.inria.fr/INRIA/publication/RR/RR-3483.ps.gz`.

[10] François Pottier. `Wallace`: an efficient implementation of type inference with subtyping, February 2000. URL: `http://pauillac.inria.fr/~fpottier/wallace/`.

[11] François Pottier. A versatile constraint-based type inference system. *Nordic Journal of Computing*, 7(4):312–347, November 2000. URL: `http://pauillac.inria.fr/~fpottier/publis/fpottier-njc-2000.ps.gz`.

[12] François Pottier. Simplifying subtyping constraints: a theory. *Information and Computation*, 170(2):153–183, November 2001. URL: `http://pauillac.inria.fr/~fpottier/publis/fpottier-ic01.ps.gz`.

[13] Didier Rémy. Projective ML. In *ACM Symposium on Lisp and Functional Programming (LFP)*, pages 66–75, 1992. URL: `ftp://ftp.inria.fr/INRIA/Projects/cristal/Didier.Remy/lfp92.ps.gz`.

[14] Didier Rémy. Type inference for records in a natural extension of ML. In Carl A. Gunter and John C. Mitchell, editors, *Theoretical Aspects Of Object-Oriented Programming. Types, Semantics and Language Design*. MIT Press, 1994. URL: `ftp://ftp.inria.fr/INRIA/Projects/cristal/Didier.Remy/taoop1.ps.gz`.

[15] John C. Reynolds. Automatic computation of data set definitions. In A. J. H. Morrell, editor, *Information Processing 68*, volume 1, pages 456–461. North Holland, 1969.

[16] Mitchell Wand. Type inference for objects with instance variables and inheritance. In Carl A. Gunter and John C. Mitchell, editors, *Theoretical Aspects of Object-Oriented Programming: Types, Semantics, and Language Design*, pages 97–120. MIT Press, 1994. URL: `ftp://ftp.ccs.neu.edu/pub/people/wand/papers/gunter-mitchell-94.dvi`.